
ONLINE BANKING FRAUD MITIGATION

*A Quantitative Study for Extracting Intelligence about Target Selection by Cybercriminals from Zeus
Financial Malware Files*

Samaneh Tajalizadehkhoob- 4181301
Engineering and Policy Analysis (ICT Management Specialization)
Faculty of Technology, Policy and Management
Delft University of Technology



THESIS COMMITTEE

CHAIR: MICHEL VAN EETEN (PROFESSOR OF PUBLIC ADMINISTRATION, SECTION POLG)

FIRST SUPERVISOR: MARTIJN GROENLEER (ASSISTANT PROFESSOR OF PUBLIC ADMINISTRATION, SECTION POLG)

SECOND SUPERVISOR: JAN VAN DEN BERG (PROFESSOR OF INFORMATION & COMMUNICATION TECHNOLOGY, SECTION ICT)

EXTERNAL SUPERVISOR: MICHAEL SANDEE (PRINCIPAL SECURITY EXPERT, FOX-IT)

Daily SUPERVISOR: HADI ASGHARI (PHD CANDIDATE, SECTION POLG)

ACKNOWLEDGMENTS

I would not have been able to complete this dissertation without the help and support of the kind people around me. Above all, I would like to thank my parents, brother and sister who not only provided me the opportunity to start with my higher education studies at Delft University of Technology, but also, support me all the way with their continuous encouragements.

The idea of this thesis topic has primarily been brought up by Prof. Michel van Eeten. From early moments in my thesis project, I have found myself lucky for working under his supervision. Thanks to Michel in both personal and academic level for his very positive role who provided me constant support and for the help in many ways throughout my project while giving me the opportunity to work in my own way.

The good advice and guidance of my first and second supervisors Dr. Martijn Groenleer and Prof. Jan van den Berg has been valuable on academic level, for which I am deeply grateful. The thesis would not have been possible without the help, patience and friendship of Ir. Hadi Aghari as my daily supervisor. I attribute the level of my master thesis to his constant encouragements, great ideas, and his friendly and truly daily supervision.

I would like to thank my dear supervisors at Fox-IT, Michael Sandee and Maurits Lucas. Thanks to them for providing me the opportunity of having access to the (or a) Zeus malware dataset which I have used as the foundation of my analyses. I was most grateful of having their great expertise and continuous help throughout my project and during my internship at Fox-IT.

I would also like to thank my dear friend Ir. Siamak Hajizadeh for his kind personal and academic help in many moments during my project especially during the process of learning the Python programming language. Special thanks to Ir. Armin Parnia for his kind and quick help on finding out extra information concerning the targeted domains in my dataset. Likewise, I would like to thank Ir. Rolf van Wegberg for his useful guidance in regards to the field of criminology.

Last, but not least, I thank all my caring friends and companions Dena, Delaram, Leila, Armin, Siavash, Gert and some of whom have already been named, whose help sustained me with hope and passion during my studies in the Netherlands.

Delft, August 13, 2013

EXECUTIVE SUMMARY

The thesis is about online banking fraud. Online banking fraud is about committing fraud or theft using the characteristics of Internet to illegally remove money from, or transfer it to, a different bank account. In general, online banking fraud is executed for the ultimate goal of gaining access to the user's bank account. How this access is obtained differs between different attack vectors. In some cases, cybercriminals trade a users' banking credentials such as PIN, password, certificates etc. for anywhere around \$10-\$2000 per account.

Until now, several technical security improvements have been achieved for securing the online banking sessions. Most of the research in this filed study the vulnerabilities that exist in the defense systems and others study the ways to increase user awareness in regards to online banking fraud. However, in this research we investigated the problem more proactively by exploring the reasons why certain banks are selected as targets for online attack and some banks are not.

It should be noted that target selection is as a socio-economic problem as technical. Reviewing economics of information security literature, we realized that some targets may be selected due to the herding behavior among cybercriminals that copycat each other's target list. To gain more systematic insights about the situations in which an online banking attack may occur and about the properties of a suitable target, we studied Routine Activity Theory (RAT) from criminology literature. As for the empirical data of our investigations, we used our unique opportunity of having access to Zeus financial malware files provided by Fox-IT Security Company.

Thus, our research question was born as *"can we extract intelligence on criminal attack patterns and target selection from the files which financial malware use as instructions for their operations?"* To answer this research question, we built a conceptual framework based on RAT. As the starting point of our research, the raw data provided by Fox-IT needed to be aggregated. In the first step of data preparation, we extracted a number of preliminary information points from the Zeus malware data files: domains targeted, active botnets (RC4keys that send configuration¹ file) and infection times (the first time that a configuration file is seen).

As for the core work in this research, by exploring the primary information from the Zeus malware data files, we identified eleven variables extractable from Zeus malware files that can provide intelligence about target selection either directly or indirectly through other variables. For each variable we have discussed the intelligence that it is able to provide.

In the process of looking at the variables, some interesting empirical patterns were found:

- ❖ Our dataset confirms the 80-20 power law distribution: a smaller group of domains (15%) attract almost 90% of the attacks.
- ❖ There are a small number of domains (88 domains- '**always-attacked**' group) that are attacked over the whole time span of our data from Jan 2009 to March 2013.
- ❖ There is a large group of domains (1170 domains- '**rarely-attacked**') that has been attacked less than 10 times in the whole period from Jan 2009 to March 2013. This implies huge target availability for cybercriminals. It could be attributed to huge amounts of trial and error and R&D by the cybercriminals. Or, the attacks are executed very selectively for specific social/political/personal reasons.
- ❖ The number of active botnets had a decreasing trend from mid-2011 onwards, which might be attributed to the Zeus take down efforts being done by national governments and security firms.

¹ A configuration (config) file gives instructions to the malware. The specifications of Zeus malware configuration file are explained in in chapter four of this report.

- ❖ Botnets differ in their attack strategy in regards to the geographical location of the target. The data indicated that among the 10 most active botnets, some attacked a wide variety of countries while others attacked more concentrated and only to a small number of countries with similar attributes

We argued that current interpretations extracted from configuration files of financial malware in regards to targeted domains in different online banking attack are incomplete. Accordingly, we intend to find a dependent variable to represent number of times a target is attacked more accurately. Therefore, we selected the variable “**average number of botnets attacking a domain**” among all of the variables introduced to be our dependent variable that is able to address target selection by cybercriminals. This variable is selected as a method for extracting and counting attacked domains from malware configuration files primarily because it was in line with the scope of our research and secondly because using this variable, the limitations of over/under counting would be mediated.

In addition, we identified a number of explanatory factors through interviewing experts of the field that may influence cybercriminals’ decisions in selecting their targets. As for the ‘proof of concept’, a series of six hypotheses were built based the empirical model presented in chapter six of this report, were statistically tested in this research. Among them, we noticed that banks with the English option in their banking webpages are attacked more within the EU region. Also countries with higher rate of GDP and broadband penetration were attacked more globally. The result of our *decision tree analysis* also confirms this.

Looking back at the research question, we can state that yes it is possible to find intelligence on criminal attack patterns and target selection from the Zeus financial malware configuration files. One of the variables being extracted from the files is statistically test in this research as the dependent variable. However, several other variables are introduced that can study target selection of cybercriminals from different aspects and thus could be interesting topic for future research.

Two points of discussion can be raised based on this research. As the results of our analysis indicate, among all of the relations that are tested, many of them did not hold. Primarily, this might imply that attack decisions of cybercriminals are dynamic and complex to the extent that even for the people in the front line it is hard to see the patterns in the attacks by criminals. Also the results can potentially questions that extent of which RAT is exposable to virtual space problems.

Table of Contents

Acknowledgments	2
Executive Summary.....	3
List of Figures	8
List of Tables	11
Chapter 1 - Introduction.....	12
1.1 Research Problem	12
1.2 State-of-the-Art Literature	13
1.2.1 Economics of Information Security.....	13
1.2.2 Criminology Theories: RAT.....	13
1.3 A Unique Opportunity for Research.....	14
1.4 Research Question	14
1.5 Outline of the Thesis	14
Chapter 2 - A Study of the Literature	17
2.1 Online Banking Today: A Brief History	17
2.2 Online Banking	18
2.3 Cybercrime and the Problem of Online Banking Fraud.....	18
2.4 Actors of Online Banking Fraud.....	20
2.4.1 Malicious Exploiters	20
2.4.2 Money Mules	20
2.4.3 Victims	20
2.4.4 Security Guardians.....	21
2.5 The Most Common Online Banking Attack Vectors	21
2.5.1 Server Side attacks.....	22
2.5.2 Client Side Attacks	22
2.6 Review of Economics of Information Security Literature	26
2.6.1 Online Banking Fraud as a Wicked Problem	26
2.6.2 Herding Behavior in Target Selection in Online Banking Fraud	29
2.6.3 Concluding Remarks	30
2.7 Review of Criminology Literature.....	31
2.7.1 Routine Activity Theory	31
2.7.2 Is RAT Applicable to Cyberspace?	31
2.7.3 RAT and Online Banking Fraud	32
2.8 Building A Conceptual Framework	34
2.8.1 A Research Gap.....	34
2.8.2 The Research’s Unique Opportunity.....	34
2.8.3 Formulating the Research Question	35
2.8.4 The Conceptual Framework.....	36
Chapter 3 - Methodology	40
3.1 Data Preparation.....	42

3.2	Extracting Intelligence from Zeus Malware Configuration Files	42
3.3	Building the Empirical Model	42
3.4	Data Analysis	43
Chapter 4	- Data Preparation	45
4.1	Zeus Banking Trojan Malware	45
4.1.1	Zeus Technical Facts	45
4.1.2	Citadel	49
4.1.3	Ice IX	51
4.1.4	Peer-to-Peer (P2P) Zeus	51
4.2	Overview of Fox-IT Malware Dataset	51
4.2.1	General Specifications	52
4.2.2	Data Collection Method	52
4.2.3	Data Format	52
4.3	Building the Research Dataset	55
4.3.1	Extracting the Key Values	55
4.3.2	Infection Time	61
4.3.3	RC4 Key (Botnet key)	61
4.4	Summary	62
Chapter 5	- Extracting Intelligence from Zeus Configuration Files	64
5.1	Discussions on Interpretations of Malware Configuration Files	64
5.1.1	Current Interpretations	64
5.1.2	Why Are the Current Interpretations Incomplete?	64
5.2	List of Extractable Intelligence from Zeus Data	65
5.3	Concluding Remarks	78
Chapter 6	- Building the Empirical Model	81
6.1	Interviewing Experts to Identify Independent Variables	81
6.2	The Identified Independent Variables	81
6.3	The Empirical Model	83
6.4	Proposed Hypotheses	84
Chapter 7	- Analysis	86
7.1	Descriptive Analysis	86
7.2	Population of the Targeted Domains	86
7.3	Botnet Activity Trends	88
7.4	Attack Persistence	90
7.5	Bivariate Statistical Analysis	92
7.5.1	Specifications of the Dataset	92
7.5.2	Level of Analysis	92
7.5.3	Determining location & language Information	92
7.5.4	Selecting Geographical Region	92
7.5.5	Pooled Data on All Years Versus Focusing on a Single Year	93
7.6	Test of Individual Hypothesis	93
7.6.1	Test of Normality	93

7.6.2	Comparison of Means.....	99
7.6.3	Measures of Association.....	103
7.6.4	Correlation Matrix.....	110
7.6.5	EU Versus Global.....	112
7.7	Decision Tree Analysis.....	113
7.7.1	Building the Dataset.....	113
7.7.2	Random Forest Test.....	113
7.8	Summary of Findings.....	115
Chapter 8 - Discussion and Implications.....		117
8.1	Reviewing the Findings.....	117
8.1.1	A Method for Expressing Popularity of Targets.....	117
8.1.2	Intelligence Extracted from Zeus Malware Configuration Files.....	118
8.1.3	Test of Hypotheses.....	119
8.1.4	Summary.....	120
8.2	Recommendations and Practical Implications.....	120
8.2.1	Scientific Recommendations.....	120
8.2.2	Recommendations for Fox-IT.....	121
8.3	Discussion.....	122
8.4	Limitations and Suggestions for Further Research.....	123
8.4.1	Reliability of Instruments.....	123
8.4.2	Validity.....	123
8.4.3	Recommendations for Further Research.....	125
Literature List.....		126
Appendix - Data Triangulation.....		129
Comparison Between Number of Configuration Files.....		129
Comparison Between Number of Configuration Files versus Number of Command and Control (C&C) Files.....		131
Comparison Between top-20 Banks Targeted by Zeus versus SpyEye Trojan.....		131
Comparison Between Top Attacked Countries.....		132
Comparison Between Proportion of Emails Contain URL Malware and Number of Zeus Configuration Files.....		133

LIST OF FIGURES

Figure 1- Routine Activity Theory (OFFENDER+ TARGET- GUARDIAN =CRIME)	13
Figure 2- online banking general components (Adopted from Hutchinson & Warren, 2003)	18
Figure 3- Online fraud actors (ADOPTED FROM FBI, 2012)	21
Figure 4- Phishing targets by industry and country (McAfee, 2012)	23
Figure 5- Example of Phishing Email contains disguised link.....	23
Figure 6- Man-in-The-Browser attack operation (Utakrit, 2009).....	24
Figure 7- Routine Activity Theory (offender+ target- guardian =crime)	31
Figure 8- Online Banking Fraud using RAT framework	33
Figure 9- Initial conceptual framework.....	37
Figure 10- Overview of the research steps.....	41
Figure 11- Zeus Infectioini cycle (TRENDMICRO, 2010)	47
Figure 12- Example of a webpage after Zeus web injection (Macdonald, 2011).....	48
Figure 13- Zeus variants' distribution (TrendMicro, 2010).....	49
Figure 14- Citadel distribution in Europe (SHERSTOBITOFF, 2013)	50
Figure 15- Worldwide Citadel distribution (Sherstobitoff, 2013).....	50
Figure 16- The P2P C&C paradigm has shifted from central server to botnet (Symantec, 2012)	51
Figure 17- The process of aggregating the Zeus malware configuration files.....	53
Figure 18- tables in the MYSQL Zeus database (red circles indicates the variables that are used for this research from each table)	54
Figure 19- Preview of Zeus malware configuration file	55
Figure 20- Overview of the main tables used in the research.....	62
Figure 21- Example for displaying why 'raw counts' are incomplete.....	65
Figure 22- Intermediate framework based on the relations in the conceptual ramework.....	66
Figure 23- Top 10 domains attacked by Zeus Malware Jan 2009-March 2013	67
Figure 24- Top 10 most targeted domains attacked by Zeus malware Jan 2009-March 2013.....	67
Figure 25- Power law (cumulative percentage of domains count for percentage of attacks).....	68
Figure 26- number of weeks that domains were under attack.....	69
Figure 27- Top 10 domains attacked by Zeus malware Jan 2009-March 2013	70
Figure 28- Top 10 most targeted domains attacked by Zeus malware Jan 2009-March 2013.....	71
Figure 29- Power law distribution (cumulative percentage of domains count for percentage of attacks)	71
Figure 30- Number of active botnets per week.....	72
Figure 31- Nu Number of config files vs. number of active botnets per week	73
Figure 32- Average number of configs per botnet per week	73
Figure 33- Number of domains attacked by all botnets per week.....	74
Figure 34 Number of domains sent per botnet per week vs. number of botnets active	75
Figure 35- Botnets lifespan vs. activity (In terms OF NUMBER of config files they sent)	76
Figure 36- The empirical model	83
Figure 37-Top-10 attacked domains by Zeus financial malware	87
Figure 38- Top-10 attacked domains by Zeus financial malware	87
Figure 39- Distribution of the Top-10 attacked domains by Zeus financial malware	88
Figure 40- Number of active botnets per week.....	88
Figure 41- Number of config files vs number of active botnets per week.....	89
Figure 42- Average number of config files per botnet per week.....	89
Figure 43- Average number of domains attacked per botnet per week	90

Figure 44- Amount of weeks domains were under attack	90
Figure 45- Alexa ranks vs. Zeus ranks	91
Figure 46- Result of test of normality	93
Figure 47- Histogram of variable 'ATTACKS'	94
Figure 48- Descriptives of variable 'ATTACKS'	94
Figure 49- Result of KS test of normality.....	94
Figure 50- Histogram of variable 'SUM OF ATTACKS'	95
Figure 51- Descriptives of variable 'SUM OF ATTACKS'	95
Figure 52- Result of the comparison of means test for H1	99
Figure 53- Result of the comparison of means test FOR H1	99
Figure 54- Result of the comparison of means test FOR H1	100
Figure 55- Result of the comparison of means test FOR H1	100
Figure 56- Result of the comparison of means test for H2	101
Figure 57- Box Plot of som of attacks within two groups of with/without English webpages	101
Figure 58- Result of the comparison of means test for H2	102
Figure 59- Result of the comparison of means test for H2	102
Figure 60- Box Plot of som of attacks within two groups of with/without English webpages	102
Figure 61- Result of the comparison of means test for H2	103
Figure 62- Result of Spearman's bivariate correlation test for H3.....	103
Figure 63- Result of Spearman's bivariate correlation test for H3.....	104
Figure 64- Scatter plot of broadband users versus sum of attacks on different countries	104
Figure 65- Result of Spearman's bivariate correlation test for H3.....	104
Figure 66- Result of Spearman's bivariate correlation test for H3.....	105
Figure 67- Scatter plot of broadband users versus sum of attacks on different countries	105
Figure 68- Result of Spearman's bivariate correlation test for H4.....	106
Figure 69- Result of Spearman's bivariate correlation test for H4.....	106
Figure 70-SCATTER PLOT OF GDP VERSUS SUM OF ATTACKS ON DIFFERENT COUNTRIES.....	107
Figure 71-RESULT OF SPEARMAN'S BIVARIATE CORRELATION TEST FOR H4	107
Figure 72-RESULT OF SPEARMAN'S BIVARIATE CORRELATION TEST FOR H4	107
Figure 73- Scater plot of GDP versus sum of attacks on different countries.....	108
Figure 74- Result of Spearman's bivariate correlation test for H5.....	109
Figure 75- Result of Spearman's bivariate correlation test for H5.....	109
Figure 76- Result of Spearman's bivariate correlation test for H5.....	110
Figure 77- Result of Spearman's bivariate correlation test for H5.....	110
Figure 78- Result of Spearman's correlation test for categorical independent variables	111
Figure 79- Result of Spearman's bivariate correlation test FOR NON-geographical independent variables.....	111
Figure 80- Descriptive Statistics within the EU region	112
Figure 81- Descriptive statistics for the global region	112
Figure 82 – Visual tree (Left part)	114
Figure 83 – Visual tree (Right part).....	114
Figure 84- Top-10 attacked domains (old metric)	118
Figure 85- Top-10 attacked domains (New metric)	118
Figure 86-Fox-IT (Left) versus F-SECURE (Right)	132
Figure 87-F-SECURE P2P ZEUS top attacked countries (Left) AGAINST FOX-IT ZEUS top attacked countries (Right) (F-Secure, 2012b).....	133

Figure 88-Proportion of Email traffic Containing URL malware from Symantec data (Left) againsts Fox-IT number of Zeus configuration files(Right)(Symantec, 2013).....133

LIST OF TABLES

<i>Table 1- Screen Shot from a carding website (Paget, 2010)</i>	19
<i>Table 2- Summary of online banking attack vectors</i>	25
<i>Table 3- Security incentives of players in the ICT value chain (Bauer & Van Eeten, 2009)</i>	28
<i>Table 4- List of factors that influence online banking crime</i>	38
<i>Table 5- General specifications of the research dataset</i>	52
<i>Table 6- Number of targeted domains in different years</i>	52
<i>Table 7- PRELIMINARY information extracted from the malware configuration files</i>	55
<i>Table 8-example of targeted Obfuscated URLs in the configuration FILES</i>	56
<i>Table 9- Number of weeks top-10 domains are under attack</i>	68
<i>Table 10- Information about geographical locations attacked by the top-10 most active botnets</i>	77
<i>Table 11-Independent variables extracted from Zeus configuration files based on RAT's three categories</i>	79
<i>Table 12- List of independent variables based on interview results, grouped using RAT categories</i>	82
<i>Table 13- Spearman's correlation test between Alexa ranks and Zeus ranks of attacked domains</i>	91
<i>Table 14- List of hypotheses</i>	92
<i>Table 15- List of EU countries in/out of our sample</i>	92
<i>Table 16- Attacked domains per region per year</i>	93
<i>Table 17- Summary of the findings</i>	111
<i>Table 18- Countries in our dataset</i>	113
<i>Table 19- Summary of the findings</i>	115
<i>Table 20-Limitations of the research</i>	124
<i>Table 21-Suggestions for further research</i>	125
<i>Table 22-List of the well-known available security reports and blogs</i>	129

1.1 RESEARCH PROBLEM

Around mid-1990's, banks started to offer Internet banking mainly to increase customer reach and cost-effectiveness (Jaleshgari, 1999). Electronic banking platforms act as efficient channels throughout which transactions could be done with less effort (Vrancianu & Popa, 2010). However, as these web-based, 'online banking' platforms have become popular among citizens and e-businesses that use them more and more in their daily activities; 'Online Banking Fraud' increased likewise (Alaganandam, Mittal, Singh, & Fleizach, 2007). Cybercriminals started to use the characteristics of Internet to generate very costly scams to steal banks' costumers account information and ultimately their money (Choo, 2011). In a study Anderson et al. (2012) indicate that the global financial losses of such activities are in the magnitudes of billions of euros per year. The study provides widespread details about cybercrime losses by categorizing them into direct losses, criminal revenue, indirect losses and indirect costs. Therefore, there is a collective consensus among defenders and policy makers that measures have to be taken to protect the online banking platforms from such threats (Anderson et al., 2012).

Although, a lot of different technical controls and defense measures were taken by banking sector and security firms for instance through increased real-time supervision on transactions, yet the online banking fraud exists as a big problem (Premchaiswadi, Williams, & Premchaiswadi, 2009). The problem is that defense measures are often reactive, time consuming, public and thus accessible for anyone. However, attackers are proactive; they allocate their time investigating the already existing defense knowledge and measures to find solutions to bridge these defense measures. So it is like an iterative loop; usually attackers find the weakest links to attack (the most vulnerable link), defenders fix the problems, attackers find a new and so a strong dynamic process is at play (Böhme & Moore, 2009).

One of the ways for designing mitigation strategy for online banking sector is to primarily investigate why certain banks are selected as targets for online attack and some banks are not. Different banks from all around the world become targets of financial banking attacks and their users suffer from malware attacks in the form of online phishing², keystroke-logging³ malwares, identity theft etc.; A study done by Moore and Clayton (2007) indicates that some banks are targeted much more frequently than others. However, there does not seem to be a clear pattern in which banks are attacked by financial malware. For instance, the attacked banks differ in terms of their size (market share) or number of users, their authentication mechanisms, their money transfer policy, regulation, the statistics of the countries they are located in and many other characteristics. Some reports published in 2012 claimed that result of their analysis on some malware samples showed that online banking attacks are getting more target-specific (Sherstobitoff, 2013; TrendMicro, 2013). However, we still do not know about whether or not cybercriminals indeed select their targets based on specific characteristics and what their rationale is for selecting their banking attack targets.

² Based on Wikipedia definition, "Phishing" is the act of attempting to acquire information such as usernames, passwords, and credit card details (and sometimes, indirectly, money) by masquerading as a trustworthy entity in an electronic communication (Wikipedia, 2012).

³ Keystroke logging (more often called keylogging or "keyloggers") is the action of tracking (or logging) the keys struck on a keyboard, typically in a covert manner so that the person using the keyboard is unaware that their actions are being monitored (Wikipedia, 2012).

1.2 STATE-OF-THE-ART LITERATURE

1.2.1 ECONOMICS OF INFORMATION SECURITY

To tackle the problem just has been explained, we will start by reviewing the main literature of the field, economics of information security. Internet security became an interesting topic for social scientists from 2000 onwards because not so many satisfactory results were produced by employing technical solutions alone (Haukson, 2006). According to Anderson and Moore (2006), in order to be able to implement efficient solutions in the field of Internet security, economic insights should be integrated into technical designs. Schneier (2011) discussed that “people often represent the weakest link in the security chain and are chronically responsible for the failure of security systems”. Here is one of the situations where economics comes in the financial online crime where multiple different actors are involved. Economics intends to explain the behaviors and incentives of attackers and defenders and consequently helps to build a safer Internet environment (Hammock, 2010).

Therefore, economics of information theories has been used many times by multiple authors such as Anderson et al. (2012) and Van Eeten, Bauer, Asghari, Tabatabaie, and Rand (2010) to address cybersecurity problems in general, including online banking fraud. In regards to the online banking fraud problem, economic theories are useful to explain the main goal of the attackers in an online banking attack that is economic and for the main reason of financial gains. However, yet economic theories alone are not able to provide the required information to explain why certain targets are selected by cybercriminals. For gaining more insight about ‘target selection’, we need to refer to a similar field, which actually deals with the same issue as we are dealing here in this research.

1.2.2 CRIMINOLOGY THEORIES: RAT

In general, criminology field covers several well-known theories in which situations where crime occurs are explained. ‘**Routine Activity Theory**’ (RAT) defined by Cohen and Felson (1979) is one of the most important criminology theories which explains that a crime happens when the three components ‘Attractive Target’, ‘Motivated Offender’ and ‘Lack of Capable Guardian’ meet each other in a specific location and time (Figure 1). Although the theory was developed for offline crime, it can similarly explain the occurrence of online banking attacks and the associated underlying reasons that resulted in online attacks. As it is obvious, using this theory, we might be able to address what makes a target attractive and how it is selected by cybercriminal in an online banking attack. Of course, exposing RAT into a virtual problem, requires it to be adjusted to the specific characteristics of cyberspace (such as existence of many-to-many connectivity, collapse of spatial-temporal barriers and anonymity of mutually dependent actors).



FIGURE 1- ROUTINE ACTIVITY THEORY (OFFENDER+ TARGET- GUARDIAN =CRIME)

1.3 A UNIQUE OPPORTUNITY FOR RESEARCH

In the previous sections we introduced the problem of online banking fraud and we argued that after all improvements in defense measure by defenders, they still do not have any clue how the attack target are selected by cybercriminals.

In order to obtain information about how targets are selected by criminals and in general online banking attacks, their technical specification and their previously selected targets, we need to have access to records of a financial malware as a case study for this research. In this research, a unique opportunity of having access to records (configuration files) of Zeus financial malware is existed. Zeus is one of the most important and wide-spread malware, first exposed on July 2007. It mostly works on computers using Microsoft Windows operating systems and since 2012, also on Blackberry and Android phones. The Zeus malware files (configuration files) contain all of the instructions that are provided by cybercriminals to the infected machines including target domains⁴.

The dataset is provided by Fox-IT for Delft university of Technology. Fox-IT is a security company located in the Netherlands, Delft that provides security solutions for more than 20 countries world-wide.

The opportunity of having the dataset provided by Fox-IT is unique primarily because normally finding information about a malware configuration file which may contain confidential information about different companies and financial institutions is hard for independent researchers. Moreover, the dataset covers Zeus records of wide variety of countries and a sufficient timespan (2009-2013).

1.4 RESEARCH QUESTION

According to the discussions above, our goal in this research is to gain possible insights about how targets are selected by cybercriminals in online banking fraud. We are mainly interested in seeing whether a clear pattern/intelligence can be found from the way cybercriminals select their targets. Therefore, in line with our goal we need to perform an explorative research on the available Zeus malware data provided by Fox-IT database (Kothari, 2009). To gain any insight, we need to primarily develop a method for extracting meaningful information that could explain target selection by cybercriminals from the Zeus malware data. After realizing the best variable for delineating the information available in Zeus malware data into meaningful terms, a framework can be built containing a dependent and a number of independent variables to determine what factors have the most influence on target selection by cybercriminals. The starting point for realizing this goal is to grasp the primary datasets and meanwhile to understand the dynamics of the Zeus attacks in more details. Therefore, the above research goal can be translated into the leading research question:

Research Question: *Can we extract intelligence on criminal attack patterns and target selection from the files which financial malware use as instructions for their operations?*

1.5 OUTLINE OF THE THESIS

The thesis report is structured as follows: In the second chapter, a review of the current literature on online banking fraud will be provided, and we shall reflect on the state of the art in different online banking attacks, economics of information security literature and Routine Activity Theory from criminology literature. The chapter will start with a look at the online banking history, and gradually shift its focus on online banking attacks and its

⁴ Zeus financial malware and its specifications will be explained in more details in chapter four of this report.

different types, and target selection. The chapter will end by recognizing the research gap, developing the research question and its associated sub questions and the conceptual framework.

In the third chapter, the methodology that is used to carry out the research is presented, and the sources of data used for compiling a usable dataset are introduced. The fourth chapter complements the third chapter by explaining what Zeus malware is, how it works and how the raw malware data is aggregated and prepared along with looking at a number of special issues involved in preparing the data.

Chapter five will present the methodological core of the thesis. The chapter mainly contains different methods in term of intermediate variables that are able to extract intelligence from Zeus malware data. The chapter will finish by selecting the best fitted intermediate variable capable of delineating target selection by cybercriminals in line with our research question. This variable will then be used in later chapters as a dependent variable for empirical analysis. Chapter six and seven actually act as the proof of concept we introduced in chapter five (the methodology for extracting intelligence from data). Chapter six contains a list of independent variables identified through conducting qualitative interviews with experts from security and criminology fields. These variables will later be used to develop a number empirical hypotheses and an empirical model for illustrating the relations between variables. The hypotheses would ultimately be tested, using the empirical model that is built in the last section of chapter six with dependent and independent variables defined earlier in the report. Chapter seven starts with the section where we will finally reach the stage of data analysis. Different sets of statistical instruments are employed, and the results are interpreted. The thesis concludes with chapter eight, which summarizes all of our findings and looks at their implications. The chapter ends by a brief review on limitations of the work that is done in this thesis and a number of suggestions for further research.

Chapter 2 - A Study of the Literature

Introduction

As it is shortly discussed in the introduction chapter, in this report we are going to address the problem of online banking fraud from target selection perspective. In the current chapter, primarily we will review a short history of online banking and its different components. Next will discuss what online banking fraud is and how it can be executed. The next section contains a short review on economic and criminology theories in regards to the cyber security problems. In this section, we aim to find out to what extent these theories can help us to build a conceptual framework for addressing the research problem that is introduced earlier in the chapter.

2.1 ONLINE BANKING TODAY: A BRIEF HISTORY

As technology advanced, most of our life's daily activities moved online. Online banking was not exceptional from this fashion. Around mid-1990s, Internet brought new alternatives to the financial markets and banks started to experience the potential of the Internet (Calisir & Gumussoy, 2008; Jaleshgari, 1999).

With new technological developments, banks started to provide online banking services that enabled customers to get connected to the bank's computer systems via Internet connections through a browser or specific application (Claessens, Dem, De Cock, Preneel, & Vandewalle, 2002). Offering online banking channels not only provides an opportunity for customers to have easy access to their banking activities but also creates cost-effectiveness for financial institutions (Claessens et al., 2002; Jaleshgari, 1999). By now, performing electronic banking activities via mobile phones through Internet access is also possible for bank's customers.

Nowadays, almost every bank provides its clients with access to their accounts over the Internet. Banks provide a different range of financial services through their Internet banking channels. Different financial Internet banking applications mostly contain money transferring services, investment services (stock, bond, and mutual funds) and currency exchange services. According to the Amit and Zott (2001) definition of value creation, Internet banking services bring efficiency mostly in terms of convenience, cost-effectiveness, functionality, speed, 24/7 availability, while requiring less staff and fewer physical branches than other customer-contact channels. Since the banking industry is highly competitive, cost management is not a luxury but a necessity for financial institutions. Study of Nevens (1999) indicates that a bank's transaction costs can drop 80% or more when handled electronically.

However, as new technologies upset traditional power balances and so does the Internet. The Internet empowers everyone including cybercriminals. Ten years ago, hackers hacked systems only to satisfy their curiosity and to gain fame and so no damage was involved. However, advancement of technology and rapid progression of the hackers' ability to access various users' systems maliciously altered their motivations from curiosity to financial motives (Alaganandam et al., 2007).

Among all, online banking platforms were not exceptional and like other online-based services have become the target of various online attacks. Although banks increased their cost-effectiveness considerably by moving their consumer and business operations to the Internet environment, online banking platforms created a new risk profile for the banking section. Cybercriminals use the characteristics of Internet to perform online scams on online banking transactions. According to the Anderson et al. (2012) study on '*measuring the cost of cybercrime*', the annual global financial losses of financial fraud activities are in the magnitudes of billions of euro.

Clearly, banks desire their customers to continue using online banking. However, whether online banking sessions are safe remains one of the biggest issues that the banking sector and security firms face. If customers would not be sure that the online banking sessions are secure, they could not trust the online banking system and might

move back to the traditional way of banking which ultimately would impose higher costs on banks. Thus, banks need to provide reasonable assurance for their customers to use online banking sessions. Before going more into the details about online banking fraud, we first need to understand the basic components of an online banking session. The next section will elaborate more on the definition of online banking and its associated components.

2.2 ONLINE BANKING

Primarily, we need to explicitly know about “online banking” and the different sides that are involved in it. An online banking platform consists of the following three main components:

1. Bank (server side)
2. Network infrastructure (Internet)
3. Bank users (client side which connect to banks via network infrastructure and via personal computers, mobile phone, tablet etc.)

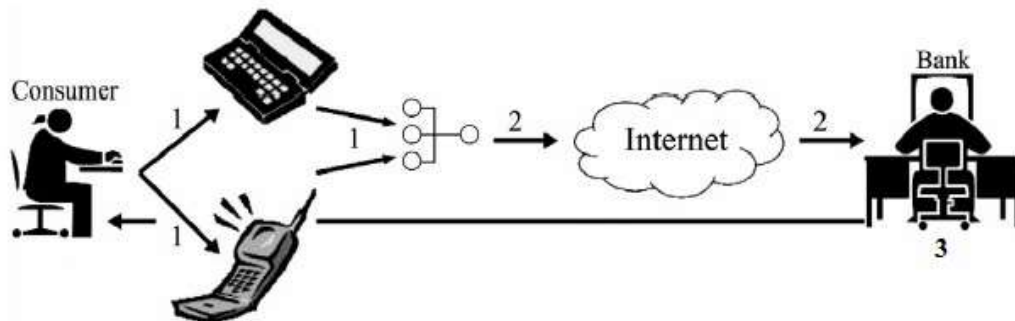


FIGURE 2- ONLINE BANKING GENERAL COMPONENTS (ADOPTED FROM HUTCHINSON & WARREN, 2003)

A bank’s user connects to the Internet typically via her personal computer, and then starts communicating with her bank via a web browser, hence banking through the World Wide Web. The standard protocol for secure online communications, HTTPS is then used (HTTP with an extra security layer).

After connecting to the main page of the bank’s online banking system, the user has to go through an authentication process to gain access to her account details and perform online any kind of transactions. Authentication methods and their level of security can differ from static (fixed) username and passwords to two step authentication which requires a token, a password (or pin) and a human being to use the token and enter the password (Claessens et al., 2002).

Indeed, each of the above-mentioned components could be potential targets for attacks by cybercriminals. However, it is logical that criminals choose the easiest way to perform an online attack and people appear to be the weakest links in this online banking value chain, since many of them do not have the required technical knowledge to protect themselves from online fraud. However, what is online fraud and how is it executed via the different components of an online banking session? In the next section, we will elaborate on these issues in more details.

2.3 CYBERCRIME AND THE PROBLEM OF ONLINE BANKING FRAUD

Douglas and Loader (2000) define, **cybercrime** as “computer-mediated activities conducted through global electronic networks which are either illegal or considered illicit by certain parties”. **Online Banking Fraud** is about committing fraud or theft using online technology to illegally remove money from, or transfer it to, a different

bank account. Wall (2001) divides cybercrime into four different categories: cyber-trespass, cyber-deceptions and thefts, cyber-pornography and cyber-violence. Online banking fraud is best fitted in the cyber-deceptions category defined as “stealing (money, property), e.g. credit card fraud, intellectual property violations (a.k.a. ‘piracy’)”. Anderson et al. (2012) differentiate online banking fraud from *card fraud* while both target financial systems and banks. They argue that in online banking fraud only customers and banks suffer while this is different in the case of card fraud where merchants also suffer from the fraudulent activity.

In general, online banking fraud is executed for the ultimate goal of gaining access to the user’s bank account. How this access is obtained differs between different attack vectors. In some cases, cybercriminals trade users’ banking credentials such as PIN, password, certificates etc. for anywhere between \$10-\$2000 per account. In other cases, the goal of the cybercriminal is to steal the victim’s money and transfer it through so called money mule accounts. Whereas other times, online banking attacks are not just about money but about harming a bank’s image by making the bank server unavailable to the real clients. An example of the first group is indicated in Table 1. The table displays a list of bank accounts for sale with different ranges of prices according to the fund available in the account.

TABLE 1- SCREEN SHOT FROM A CARDING WEBSITE (PAGET, 2010)

Bank Name	Country	Balance	Price
Bank of America	USA	---	Sold
Asmouth Bank	USA	\$16,040	€700
Washington Mutual Bank	USA	\$14,400	€600
Washington Mutual Bank	USA	\$7,950+£2,612	€500
Washington Mutual Bank	USA	---	Sold
MBNA America Bank	USA	\$22,003	€1,500
Banco Bradesco S.A.	Brazil	\$13,451	€650
Citibank	UK	£10,044	€850
NatWest	UK	£12,000	€1,000
BNP Paribas Bank	France	€30,792	€2,200
Caja de Ahorros de Galicia	Spain	€23,200	€1,200
Caja de Ahorros de Galicia	Spain	€7,846	€500
Banc Sabadell	Spain	€25,663	€1,450

Till now, several security improvements have been achieved for securing online banking sessions. Most of the research in this filed study the vulnerabilities that are existed in the defense systems (Claessens et al., 2002; Florêncio & Herley, 2011; Hutchinson & Warren, 2003; McCullagh & Caelli, 2005). Other studies focused on user part of the banking platforms and investigate the ways to increase user awareness in regards to online banking fraud. However, not so many studies, if any, look at this problem proactively by investigating the ways that target banks are selected by cybercriminals.

One of the ways for designing mitigation strategy for online banking sector is to primarily investigate why certain banks are selected as targets for online attack and some banks are not. Different banks from all around the world become targets of financial banking attacks and their users suffer from malware attacks. A study done by Moore and Clayton (2007) indicates that some banks are targeted much more frequently than others. However, there does not seem to be a clear pattern in which banks are attacked by financial malware. For instance, the attacked banks differ in terms of their size (market share) or number of users, their authentication mechanisms, their money transferring policy, regulation, the statistics of the countries they are located in and many other characteristics. Some reports published in 2012 claimed that result of their analysis on some malware samples showed that online banking attacks are getting more target-specific (Sherstobitoff, 2013; TrendMicro, 2013). However, we still do not know about whether or not cybercriminals indeed select their targets based on specific characteristics and what their rationale is for selecting their banking attack targets.

In order to be able to understand online banking fraud and the behavior of the attackers and defenders in this environment, we first need to describe the main actors involved in that space, i.e. who are the cybercriminals. The next section describes the different actors that are involved in online banking fraud value chain.

2.4 ACTORS OF ONLINE BANKING FRAUD

Actors in a malicious transaction in online banking fraud can be categorized into four main groups namely, *malicious exploiters, money mules, victims and security guardians*⁵. Below a brief description in regards to each group of actors is provided:

2.4.1 MALICIOUS EXPLOITERS

According to the anonymity that the Internet environment provides for its users, tracking people who conduct criminal activities in cyberspace is difficult. Likewise, defining a set of fixed characteristics of individuals or groups of people involved in different kinds of online fraud is a difficult task. However, According to an OECD (2007) report, based on their capabilities and motivations, malicious exploiters can be categorized into five below groups:

- **Innovators** are challenge seeking individuals who spend their time finding security holes in systems to overcome protection measures.
- **Amateur Fame Seekers** are beginners in the security world with limited programming and computing skills. Their main motivation is fame in media and they mostly use ready-made tools and tricks for executing their malicious activities.
- **Copcats** are hackers and malware authors who desire celebrity status in the criminal community and are interested in recreating simple attacks.
- **Insiders** are defined as ex-employees and ex-contractors of financial institutions or companies who are motivated to take revenge or to steal money from where they worked. Most of the times they execute their attacks by abusing their security privileges.
- **Criminals** are highly motivated, highly organized, very knowledgeable and powerful individuals or group of individuals and are playing the game for profit using advanced human and computer resources.

2.4.2 MONEY MULES

By definition of OECD (2007), **Money Mules** are “individuals recruited wittingly and often unwittingly by criminals, to facilitate illegal funds transfers from bank accounts”. FBI defined a money mule network as a network of compromised individuals who engage in transferring the stolen money from the victim’s account in exchange of gaining a percentage of that money. According to a Fox-IT security expert, in some of the reported cases, money mules are young high school boys who have been paid for lending their debit card for a certain amount time. Florêncio and Herley (2010) mentioned that the key role of money mules is to convert “reversible traceable transactions into irreversible untraceable ones”.

2.4.3 VICTIMS

Based on the OECD (2007) report victims can be categorized into two main groups: End users and Banks. End users then are categorized as home users (and individuals), small to medium Enterprises (SMEs) and large end users (public institutions and global corporations).

Home and SME users create the most negative externalities among the legitimate actors. Part of the externalities is absorbed by other players, such as ISPs (Internet service providers) and FSPs (financial service providers). Risky

⁵ The value chain economy of malicious actors is a bigger than what I have explained but the focus of this research is only on this part.

online behavior of these users along with them not employing security software creates negative externality⁶. Different reasons may be able to explain this behavior of users: many times users do not have the required knowledge about how malware works, and thus they may not understand when they become infected. Moreover, everyone perceives the chances of being hit by a malware as low, for example they believe that it should be already included on the computer. Finally, following the above points, they do not find paying for security software convenient. Moreover, because they are the *weakest link* in cyber security environment since most of them are not familiar with information technology and security requirements are too difficult for them to follow to protect themselves from Internet threats (Asghari, 2010; Mannan & van Oorschot, 2008). Considering that, some of the online threats like spamming and botnet activities are performed by acquiring users' device.

2.4.4 SECURITY GUARDIANS

Security guardians are also one of the main actors in the online banking fraud field. They improve the existing banking systems and their vulnerabilities and develop more sophisticated security systems for mitigating the online fraud on the financial institutions. This security guardians in this problem is either the banking sector itself or 3rd party security firms who provides security services for the banking sector.

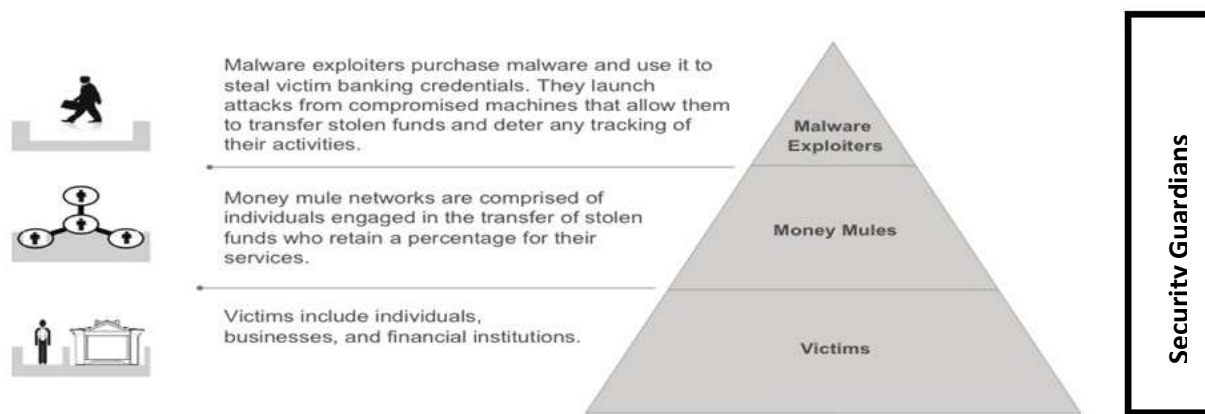


FIGURE 3- ONLINE FRAUD ACTORS (ADOPTED FROM FBI, 2012)

2.5 THE MOST COMMON ONLINE BANKING ATTACK VECTORS

In general, online banking attacks can take place using different attack methods. There is no clear-cut list of different types of online banking attacks. Primarily, because online banking attack types are not static and they evolve over time. Secondly, because categorizing different attack vectors is not an easy task to do due to the elements different attacks may have in common with each other. However, based on the targeted side, online banking attacks can be categorized into two groups: 1. Attacks that are executed on the bank's servers; 2. Attacks that are executed on the bank's clients or users.⁷

⁶ Wikipedia defines externality as a cost or benefit which results from an activity or transaction and which affects an otherwise uninvolved party who did not choose to incur that cost or benefit.

⁷ We will not include card skimming and/or attacks on merchants (third-parties) that are then used to steal credit card data in this research

2.5.1 SERVER SIDE ATTACKS

One of the types of the online banking attacks on the server side is directly hacking a bank which normally does not frequently happens because it is hard to be executed due to sophisticated banks' security infrastructure. Another type is to target bank's servers through DDoS attacks to either block the communication between the server and its real clients or to make the server unavailable to its real clients.

'Denial of Service' (DoS) or Distributed Denial of Service (DDoS) attack is to temporarily make a service or machine unavailable for its users. This is done through different ways, one of the most popular ways of DDoS attack is saturating the target service or machine through multiple anonymous service requests which ultimately result in the targeted service getting overloaded by requests. It can be a "single-source" attack, originating from only one host (attacker), or "multi-source", in which a group of hosts collaborate to flood the system with several requests (Hussain, Heidemann, & Papadopoulos, 2003). The goal of such an attack is either to interrupt the service that is provided by the target or to interrupt communication of the target with other users of the network so that it will not be able to communicate anymore. Slow network performance and speed, unavailability of some specific websites, several disconnections of one's wireless or wired network connection and receiving several spam email in a specific time are the symptoms of a DDoS attack.

This type of attack cannot easy to be executed since bank's servers are mostly well protected⁸. The fact that banks' servers are in the control of banks and more security measures are implemented for them compared to the clients' PCs make them less vulnerable targets for the attacks (although still several server side attacks are reported each year).

2.5.2 CLIENT SIDE ATTACKS

Most of the online banking attacks being executed on the client side are based on deceiving users to steal their banking credentials, PIN and ultimately their money. However, the whole process is not that easy. Several different techniques exist for each of these steps. Generally, as Anderson et al. (2012) argue in their article '*Measuring the cost of Cybercrime*', online banking fraud is carried out in two ways:

- Attacks utilizing phishing
- Attacks utilizing malware

In the following section each of the above methods will be explained in more details.

ONLINE BANKING ATTACK UTILIZING PHISHING

In a phishing attack, cybercriminals imitate a bank's website and create a fake webpage identical to the original banking page, to deceive users to provide their login credentials (Anderson et al., 2012). In their study, Moore and Clayton estimated that between 280,000 and 560,000 people gave away their credentials to phishing websites each year (Moore & Clayton, 2007).

What is a phishing attack? '**Phishing**' is a type of online threat and an example of a *social engineering* technique in which the attacker also known as the phisher attempts to fraudulently extract users' sensitive confidential information such as password, credit card information, online banking credentials etc. or to direct clients through links to download a Trojan-family on their computers for other fraudulent activities by masquerading itself as a trustworthy entity in an electronic communication (Jakobsson & Myers, 2006). Figure 4 displays the major phishing targets in 2012.

⁸ This does not imply that DDoS attacks never happen but the number of such incidents are not significant.



FIGURE 4- PHISHING TARGETS BY INDUSTRY AND COUNTRY (MCAFEE, 2012)

Usually, the electronic communication is done via spam emails or instant messaging. Phishing occurs by clicking on the links that are normally existed in the spam emails received by clients in their mailbox. The spam emails are normally distributed via bot or zombie computers in a botnet⁹. Figure 5 displays an example of a phishing email containing a disguised link that directs users to a malicious website where malware is then downloaded automatically to the user's computer.

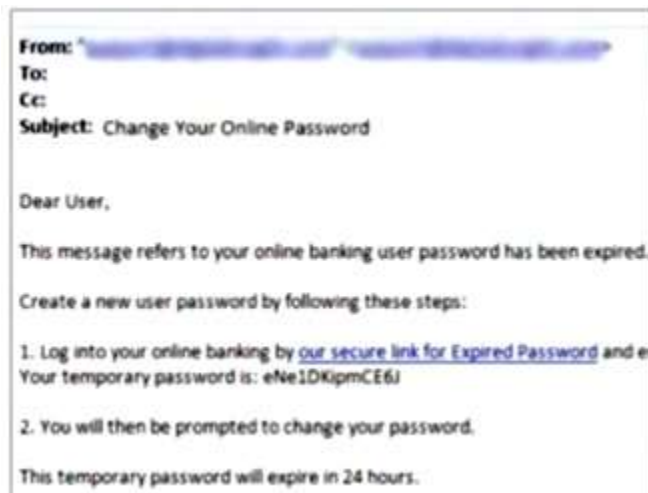


FIGURE 5- EXAMPLE OF PHISHING EMAIL CONTAINS DISGUISED LINK

Previously, in online banking fraud, phishing websites only present the bank's login webpage and it was the exact copy of the original bank's online login webpage. However, as phishing attacks are becoming more sophisticated, they rely more on malicious proxy or auto-configuration scripts that are remotely hosted by C&C servers and instruct the victim's web browser to proxy certain requests according to the specified configuration.

⁹ A 'botnet' is a flexible remote-controlled network containing computers that function together to make a platform available for fraudulent and criminal purposes (Bauer & van Eeten, 2009). In a botnet, bots are connected to the botmaster through a command and control (C&C) channel. The main value of this channel is its ability to broadcast the botmasters' commands to their bot armies (Fabian & Terzis, 2006).

ONLINE BANKING ATTACK UTILIZING KEYSTROKE-LOGGING MALWARE

Another way of performing online banking fraud is to install keystroke-logging malware on the victim's computer. Zeus malware is an example of such a malware. In October 2010, the bot master of Zeus botnet (Zbot) malware, the financial malware that harvests credentials from many banks was arrested by FBI.

Based on a McAfee report, this type of online fraud is mainly based on downloading malicious software or 'malware', typically a Trojan, on the victim's computer to control his browsers through injected scripts (Marcus & Sherstobitoff, 2012). The malware may log a user's keystroke to harvest credentials for electronic banking (Moore, Clayton, & Anderson, 2009).

The malware is normally downloaded through a link in a phishing email (spam email) that is generated by the relevant botnet (Mavrommatis, Provos, Monrose, & Abu Rajab, 2008; Moore & Clayton, 2007). The Trojan then installs the malware on the victim's device. The Next time that the victim tries to login to his/her bank account, the malware detects the victim's information and credentials. Typically, a man-in-the-browser (MitB) attack is executed by injecting scripts into the user's browser and manipulating the victim by displaying fake bank errors and a fake account balance. The malware then uses the victim's credentials to transfer money from the victim's account to a money mule account and finally forward it to a specific destination via an anonymous money transfer service, such as Web Money, e-gold, or Western Union (Paget, 2010). These mules are persons typically duped to accept stolen money and then forward it to specific destinations. The mules can easily be recruited via job ads sent in spam e-mails or hosted on websites such as Craigslist or Monster. The whole process takes 1-3 days.

What is a 'man-in-the-browser' attack? The '**man-in-the-browser**' attack also called 'Proxy Trojan', functions on Internet browser of users' computer and controls the ingoing and outgoing content of information at the system level (Utakrit, 2009). Therefore, it has the ability to take control over data or steal data either passively through key-loggers or actively through phishing. In fact, the attacker has access to whatever data that is communicated through the user's browser while he can also be selective in choosing the target domains (webpages) and the type of data he intends to steal from users in that website (Dougan & Curran, 2012).

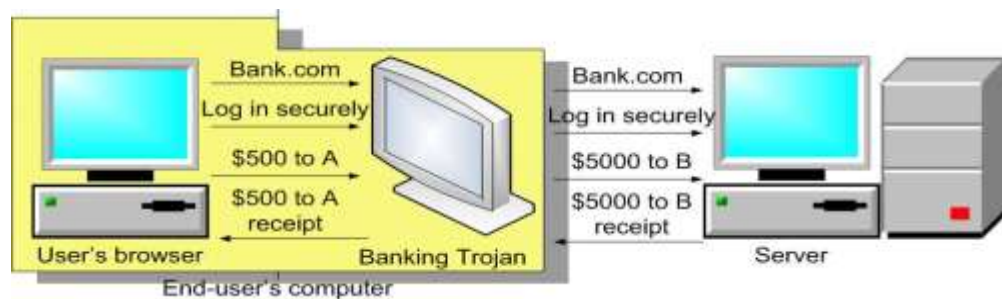


FIGURE 6- MAN-IN-THE-BROWSER ATTACK OPERATION (UTAKRIT, 2009)

Moreover, the attacker is able to inject new scripts into the user's browser and modify the HTML of specific pages to ask extra information from users. The same process could be done by the attacker for the outgoing data. This means that with MitB attacks, the attackers are able to compromise the outgoing data - the data that user submits to the server.

Table 2 provides the summary of specifications of attacks vectors explained above. Also, the associated security techniques and mechanisms that are used in each of the attack vectors are explained.

TABLE 2- SUMMARY OF ONLINE BANKING ATTACK VECTORS

Attack Goal	Attack Type	Tools and Techniques that are used	Attacked Point
Blocking the ingoing and outgoing communication of the bank's server/ harming a bank's image/ Stealing info, payment history etc.	Denial-of-Service	<ul style="list-style-type: none"> ▪ Malware ▪ Botnet 	Server
Credential/Money theft	Online banking attack utilizing phishing	<ul style="list-style-type: none"> ▪ Phishing ▪ Social engineering ▪ Spam email (containing fake website) ▪ Botnet (sending spam emails) 	Client
Credential/Money theft	Online banking attack utilizing keystroke-logging malware	<ul style="list-style-type: none"> ▪ Man-in-the-browser ▪ Proxy Trojan malware (e.g. SpyEye ,Zeus) ▪ Key-logger 	Client

In this section, we looked at online banking platforms and its different components. We discussed the ways that online banking attacks are mostly executed. We also discussed that most of the studies until now focused on the defense measures employed by the banking sector and the users themselves (users' online awareness). However, not any known study worldwide focuses on the ways criminals behave in selecting their targets. Or more specifically no studies till now investigates the patterns in which the targets are selected by cybercriminals. In order to find out the specific theories or concept that could be helpful in addressing the target selection in online banking fraud, in the next sections of the report we are going to review two relevant fields to our research namely economics of information security and criminology.

2.6 REVIEW OF ECONOMICS OF INFORMATION SECURITY LITERATURE

In the previous sections of this chapter, we discussed online banking fraud and the different attack vectors that are existed for executing online banking attacks. In this section, we will now shift to study the relevant literature to find out the explanatory factors that would possibly influence specific targets to be selected by cybercriminals. For this, we will start by reviewing Economic of Information Security literature because several times it has been argued by well-known authors of the field like Ross Anderson that information security should not be studied only as a purely technical problem but as socio-economic problem.

2.6.1 ONLINE BANKING FRAUD AS A WICKED PROBLEM

Part of the deficiencies in the decisions regarding online banking fraud problem can be attributed to the nature of this problem. The problem of online banking fraud and selection of the targets by cybercriminals is not simple but rather a complex problem. According to Rittel and Webber (1973), science is developed to deal with 'tame problems' while not all of the problems are tame but 'wicked problems'. De Bruijn and Ten Heuvelhof (2012) defined two characteristics for wicked Problems: there is no consensus about the norms regarding the problem and the information available is not objectifiable. This is true in the case of online banking where no predefined norm exists. The reason could be attributed to the facts that:

1. Online banking fraud is a new field and it influenced by several unique characteristics of cyberspace.
2. Involvement of several different actors with different interests makes this problem complex. Online banking fraud is like a game between defenders and attackers that suffers from a series of socio-economic structural asymmetries which ultimately makes the information around the game hard to objectify.

In the following section, we will review how Economic of Information Security scholars have tried to understand the complexity of online banking fraud beyond the technical explanations until now. The main goal of our review in this section is to find out the relevant factors that can influence the attack decisions of cybercriminals in selecting the target domains in online banking crime through understating the complexity of the online banking fraud problem.

REASON ONE: UNIQUE CHARACTERISTICS OF CYBERSPACE

The first reason why online banking is a wicked problem can be attributed to the unique characteristics of cyberspace. Primarily, the fact that cyberspace provide opportunity for its users to interact with spatially distant actors, apart from its novel benefits, makes users in the cyberspace more vulnerable to a range of potential cybercriminals; Considering the fact these cybercriminals have instant reach, unimpeded by the normal barriers of physical distance.

Secondly, the ability of communicating with thousands of individuals instantaneously empowers individuals with minimal resources to create potentially huge negative effects .The instances could be realized in 'spam' or 'phishing' emails and DDOS attacks (Yar, 2005). In addition, the anonymity provided by cyberspace that allows users to reinvent new social identities far from their real-world identities is another factor that makes cybersecurity problems more complex and thus harder to solve.

REASON TWO: EXISTENCE OF SOCIO-ECONOMIC STRUCTURAL ASYMMETRIES

As Anderson et al. (2012) mentioned in their paper, the security world is similar to a structured game in which some structural asymmetries are exist. We divide these asymmetries into four main groups as below: *Information, innovation, incentive and evaluation.*

Information

Anderson, Moore, Nagaraja, and Ozment (2007) argued that information in the online security space is asymmetric in the sense that defense measures are public and so accessible for anyone. However, attackers can take their time to come up with a new attack method. This causes the defenders to be reactive and attackers to be proactive. When a new attack comes up, defenders need to defend with a new defense method. Whereas, attackers allocate their time investigating the already existing defense knowledge and measures, to find solutions to breach the defense measures. Thus, the only strategy for defenders is to play catch-up and keep pace with the attackers.

On the other hand, the information is asymmetry in the market of security products as well. Akerlof (1970) attributed a market with asymmetric information to a used car market. The similarity is between buyers in used car market and users in the information security space. The car buyers do not have too much information and they cannot measure the car's quality (they cannot fully determine the quality difference between good and bad used cars), so they prefer not to pay a lot for a car. Thus, what happens at the end is that a bad car sold at the same price as a good car and eventually bad quality drives out the good quality. This is the case in the information security environment as well. When normal users cannot easily judge the quality of the antivirus they buy, they are not willing to pay a lot of money and buy expensive antivirus software. Consequently, this stimulates security software vendors to drop their costs to the market costs to sell their products rather than to invest more for secured products. This leads to existence of less actual security features in the antivirus/firewall products and ultimately increases the chance of users of those products for being potential targets that can easily be infected by different kind of malware.

Another major problem in the field of online banking security is the existence of '*Inaccurate Information*' which can ultimately results in inaccurate measures and countermeasures in the banking sector and third party security firms. The problem is that existing statistics on security failures (online baking crime, spam, phishing, number of botnets etc.) are collected by parties with incentives to under or over-report (Moore, Clayton, & Anderson, 2009). Security firms are incentivized to over-report the statistics to advertise for their products while governments have the tendency to under-report the statistics due to its possible political and social consequence.

Innovation

Creating new defense methods in the banking sector and new business models for banks as the defenders is expensive and huge sunk costs are involved to create the new infrastructures for the technology. As long as the drivers of being more innovative for the security industry exists, which mainly consists of the existence of bad guys and regulators, they would produce new technologies to overcome the security problems. However, it is not the same case for attackers. Any new trick they can come up with, capable of breaking security walls, can work for them. Thus, they are agile and they do not necessarily stick to one particular technology (Anderson et al., 2012).

As for security software vendors to obtain market dominance and to create positive network externality¹⁰, their products should be as attractive to complimentary vendors as to their customers; while vendors offering secured software often put burden on their complimentary firms. Therefore, security is neglected in the first releases and only added later when the market dominance is obtained.

¹⁰ 'Network externality', or network effects, or demand-side economies of scale, is the effect that one user of a good or service has on the value of that product to other users. This leads the value of a product or service to become dependent on the number of others using it (Shapiro & Varian, 1998). Positive network externalities occur when an increase in the number of users of a network encourages other customers to join that network or brand due to its popularity.

Gaining market dominance by having the advantage of first mover is also another fact that can explain the philosophy of “ship it Tuesday and get it right by version 3” which ultimately lead to less secure software, taking the classic example of Microsoft Windows (Moore & Anderson, 2011; Van Eeten & Bauer, 2008). Creation of such deficiencies by security software vendors gradually results in insecure an online banking session which ultimately hurts banks and their customers.

Moreover, as (Moore et al., 2009) discuss, in the Internet environment insecurity can be seen like environmental pollution, as it generates negative externalities. Every single insecure computer that is connected to Internet can be used by fraudsters to execute online fraud and therefore create negative externality.

Incentives

One of the most important structural asymmetries that resulted in the failure of security systems was first introduced by Anderson (2001). He explained that failures of information systems are due to stubborn ‘incentives’ of actors as much as it is due to technical failures. He argued that, banking systems are subject to failures because the person who guards them is not the person who suffers when they fail.

Bauer and Van Eeten (2009) categorize incentives of financial service providers into two distinct groups: **Security-Enhancing** ones versus **Security-Reducing** ones. The red box in the Table 3 illustrates these two categories for financial service providers.

TABLE 3- SECURITY INCENTIVES OF PLAYERS IN THE ICT VALUE CHAIN (BAUER & VAN EETEN, 2009)

Player	Security-enhancing	Security-reducing
Internet service providers (ISPs)	Cost of customer support Cost of abuse management Cost of blacklisting Loss of reputation, brand damage Cost of infrastructure expansion Legal provisions requiring security	Cost of security measures Cost of customer acquisition Legal provisions that shield ISPs
Software vendors	Cost of vulnerability patching Loss of reputation, brand damage	Cost of software development and testing (time to market) Benefits of functionality Benefits of compatibility Benefits of user discretion Licensing agreements with hold-harmless clauses
E-commerce providers (banks)	Benefits of online transaction growth Trust in online transactions Loss of reputation, brand damage	Cost of security measures Benefits of usability of the service
Users	Awareness of security risks, realistic self-efficacy, exposure to cybercrime	Poor understanding of risks, overconfidence, cost of security products and services

However, in case of banking section, on the security side, the incentives are naturally aligned. Nowadays, financial intuitions are one of the sectors in the security field that have their incentives in line with those of customers due to the economic benefits that online banking sessions brought for them. Bauer and Van Eeten (2009) argued that because implementing security measures and creating trust in online banking platforms would put a burden on banks’ users, banks have to make a trade-off between offering a safer online transaction platform and benefiting from a having user-friendly and convenient platform. However, based on the Fox-IT experts’ opinion, in today online banking environment most of the banks around the world, specially, in North America and Western Europe provide the optimal level of security¹¹ while at the same time they internalized the cost of online fraud for individual customers.

¹¹ By optimal level of security we mean the balance between security cost and benefits of security (e.g. Usability Vs. Security)

The rationale behind it is that, because users are the weakest link in the chain of online banking, they may not be able to assert the security of their online session themselves. Therefore, by providing the optimal level of security in the online banking sessions and refunding their loss in cases of fraud, they do their best in gaining the customer's trust. Whereas, business customers can handle security measures themselves due to the facilities, the knowledge and expertise they have. So, unlike many other security fields in which misaligned incentives of legitimate actors may amplify security failures, in this field it acts as a counter-force and helps the security to be improved to a certain extent.

It should be noted that the incentives of actors in online banking fraud game are not aligned when it comes to providing accurate information about banks being targets of online attacks which ultimately poses what Olson (2009) called it '**Collective Action Problem**'. The problem occurs because most of the times when incentives of individual businesses, governmental companies and security firms are not aligned. Individual companies do not have any incentive neither to report that they were attacked nor to share their security information with others. This is mostly because many companies are more concerned about the impact of such disclosures on their image and market value rather than on the benefit of sharing such information.

Whereas the fact is that an insecure infrastructure is bad for everyone, so sharing the information about any insecurity would be in the collective interest. An effective disclosure of security information would increase cyberspace security and general efficiency, by transferring the knowledge and experience to others that may face the same problem. However, the opponents believe that widely disclosing security information and would give the cybercriminals the advantage of having access to all the security bugs and vulnerabilities (NYTimes, 2013)¹². For instance, concerning the problem of online banking as Moore and Anderson (2011) mention, it is hard to tell which banks are suffering from any given mule-recruitment campaign. Therefore, that can explain why very few information exists in regards to attacks to different financial institutions, which eventually makes it hard for scientists of the field to determine when and where certain targets are attacked without any access to empirical attack records.

Evaluation

It is hard to evaluate how well defense measures work because few security metrics exist for measuring how secure software is. Therefore, companies in the security industry compete on every aspect of security rather than security of software itself. This is while attackers have a great way to measure the cost and benefit of the attacks they intend to execute and it is the attack monetization that drives attack quality.

The way that cybercriminals evaluate their attacks can actually explain certain factors that may influence some targets to be selected in online banking attacks by cybercriminals. For instance, economics is able to explain why banks with higher number of users are attacked more than others; simply because the benefits of attacking platforms with large number of users is higher than cost of executing such attacks (Anderson et al., 2012).

2.6.2 HERDING BEHAVIOR IN TARGET SELECTION IN ONLINE BANKING FRAUD

We should mention that many of the actions that happen in cyberspace by different actors could be a result of socio-economic deficiencies. However, mentioning these deficiencies would enhance the analysis of an individual's behavior, in this case cybercriminals. One of these theories is called "herding" in the social science literature.

¹² <http://www.nytimes.com/roomfordebate/2013/02/21/should-companies-tell-us-when-they-get-hacked/countering-security-breaches-is-a-national-priority-and-a-business-priority>

Raafat, Chater, and Frith (2009) define, herding as “a form of convergent social behavior that can be broadly defined as the alignment of the thoughts or behaviors of individuals in a group (herd) through local interaction and without centralized coordination.” In other words, herding is described as an imitation behavior resulting from individual factors which often leads to inefficient outcomes for the market as a whole. As Raafat et al. (2009) stated, the definition is generic and can be applicable to other fields as well. For instance, in regards to online banking fraud, herding behavior among cybercriminals can cause a bank or set of banks to be selected as a target for attack by criminals only because it has been a selected as before a target in previous attacks by other group of criminals. This is true taking into account the fact that attackers often speak about their successful attacks in underground forums or chatrooms.

This behavior of cybercriminals can be more specifically explained by the concept of “**informational cascade**”. Bikhchandani, Hirshleifer, and Welch (1992) discussed in their article that information or informational cascade occurs when an individual imitates the behavior of other(s) independent of his own information only because the decision or behavior other(s) seems more optimal. This means that as soon as an individual ignores his own information in favor of others’ earlier decisions or actions, the cascade develops. This can cause a very rapid trend in which effect of behavior of an individual is amplified only because it has been imitated by a mass of individuals that behave based on earlier decisions or information of their mates.

To explain the relevance of such a theory within our research, one should note that one of the objectives of this research is trying to find the independent factors that would explain patterns in selection of targets by cybercriminals. The factors later will form hypothesis, which will be the basis of our analysis. Considering the “information cascade” theory, we can argue that, a number of attacks to a target or targets have happened only because cybercriminals copy the targets other group of cybercriminals that executed successful attacks earlier in time of but not according to the bank-specific or country-specific characteristics.

2.6.3 CONCLUDING REMARKS

In the above section we made a short review of economics of information security literature in order to find a theory capable of explaining incentives of cybercriminals in selecting their target selection in online banking attacks. However, economics literature was only able to provide fragmented observations about why certain targets are more attractive such as herding behaviour among cybercriminals and about underlying economic incentives of cybercriminals in selecting their targets.

However, what we need is a more systematic way to identify different factors that would be able to explain the situations where attack decisions of cybercriminals in selecting their targets are actually take place. That is why in the next section of the report we shall review the criminology literature to see whether it is possible to extract relevant theories in regards to target selection. Specifically, we chose to review criminology literature because we find it a relevant field capable of explaining the situations in which a crime actually happens in the non-virtual world.

2.7 REVIEW OF CRIMINOLOGY LITERATURE

2.7.1 ROUTINE ACTIVITY THEORY

Routine Activity Theory (RAT) defined by Cohen and Felson (1979) is a general theory of crime causation focused on situational crime. The theory mainly explains that for crime to be committed, three aspects are needed: 1) a **Motivated Offender** 2) a **Suitable Target** and the 3) **Absence of Capable Guardian**. Cohen and Felson (1979) emphasized that the ability of the above elements to explain patterns of offending depends on these elements converging in time and space. In other words, as Yar (2005) also indicates, “routine activities always occur in particular locations at particular times”.

The theory hypothesizes that “criminal acts require the convergence in space and time of likely offenders, suitable targets and the absence of capable guardians” (Cohen & Felson, 1979). To what extent routine activity theory can embrace novelties within cyberspace is dependent on the match between structural characteristics of the virtual and non-virtual world. In the next section we will elaborate more on the extent that RAT can be applied to cyberspace.



FIGURE 7- ROUTINE ACTIVITY THEORY (OFFENDER+ TARGET- GUARDIAN =CRIME)

2.7.2 IS RAT APPLICABLE TO CYBERSPACE?

The section will start by recalling the definition of cybercrime from first part of this chapter. As it is already discussed, cybercrime is defined as “computer-mediated activities which are either illegal or considered illicit by certain parties and which can be conducted through global electronic networks” (Douglas & Loader, 2000). The trait of cybercrime is therefore inherent in the unique interactional environment in which it takes place, namely the ‘virtual space’ or ‘cyberspace’ generated by the interconnection of computers into a worldwide network of information exchange (Castells, 2003; Yar, 2005).

The differences between a number of structural properties of the virtual and non-virtual worlds, limit the extent in which RAT can be applicable to cybercrime. The difference can be mainly attributed to cyberspace’s novel socio-interactional characteristics discussed by Yar (2005) such as the collapse of spatial–temporal barriers, many-to-many connectivity, and the anonymity and plasticity of online identity.

To shortly address the differences in more details, we will discuss the application of the three criteria of RAT in cyberspace. According to Yar (2005), the first criterion, ‘*motivated offenders*’, can to a large extent be assumed

homologous in the virtual and non-virtual world and their motivation is assumed to be given. *'Suitable Target'* criterion is evaluated by cybercriminals according to four-fold essential properties namely, value, inertia, visibility and accessibility.

In terms of *'value'* property, the virtual and non-virtual worlds are similar in the sense that targets will vary due to the shifting values devoted socially and economically to specific goods at specific times. Certainly, factors such as scarcity and fashion will play a role in setting the value of a target by offenders as well as others (Cohen & Felson, 1979). The same can hold true for the virtual world where targets are valued by the reward they can bring in.

'Inertia' or portability in non-virtual world is defined as the physical properties of objects or persons that might offer varying degrees of resistance to effective predation. Yar (2005) argued that *'Inertia'* in cyberspace can firstly be addressed to the volume of data (e.g. file size) that influences the portability of a target and secondly to technical specifications of target (e.g. hard-drive space in case of information theft) which may influence the associability of target. For instance in terms of online banking attack it can be attributed to the cross-national money transferring policy of the target banks.

In terms of *'visibility'* in traditional crime, Yar (2005) argued that "property and persons that are more visible are more likely to become targets". Whereas in the case of cyberspace, physical distance is meaningless and virtually present entities can always be visible to a pool of motivated offenders.

'Accessibility', defined as the "ability of an offender to get to the target and then get away from the scene of a crime" by Cohen and Felson (1979), is not correlated to cyberspace because physical accessibility does not matter. However, a correlation exists between the two when it comes to addressing the security devices or mechanisms such as encryption devices and passwords, etc., that prevent unauthorized access to cyberspace elements.

Lastly, *'capable guardians'* in the non-virtual world is grouped to social-informal guardians such as neighbors and homeowners, or formal ones such as the police and army. The virtual world is also characterized by these capable guardians in both categories: social-informal guardians like in-house network administrators, systems security staff, ordinary online citizens and physical or technological guardians such as firewalls, intrusion detection systems and virus scanning software (Denning, 2000).

In conclusion, the differences in application of Routine Activity Theory (RAT) between the virtual and non-virtual worlds require that RAT's concepts be adapted for use in cyberspace. This implies that, while analyzing a cybercrime, one should always take into account the so called *'novel'* characteristics of cyberspace to be precise, such as failure of spatial-temporal barriers, many-to-many connectivity, and the anonymity and plasticity of online identity (Yar, 2005). In the next section, we will see how online banking fraud can be delineated using RAT.

2.7.3 RAT AND ONLINE BANKING FRAUD

In the previous section, we concluded that RAT could be exposed to the cyberspace in cases being customized and adapted with cyberspace characteristics. Accordingly, RAT can be one of the basis theories used in this research. In this section, we will apply RAT into the concept of target selection by cybercriminals in online banking fraud.

In online banking fraud, a *'Suitable Target'* is the bank's websites or other type of web pages that are included in the malware configuration files to be targeted in online banking attacks. *'Motivated Offender'* in this case is the cybercriminal who plans and executes the attack. The *'Absence of Capable Guardian'* is enhanced by cyberspace's anonymity characteristic and could be defined as the absence of security counter measures by banks and the 3rd parties who provide those measures such as Fox-IT. Accordingly, online banking fraud (OBF) takes place where these three criteria meet. Figure 8 displays the occurrence of online banking fraud using RAT.

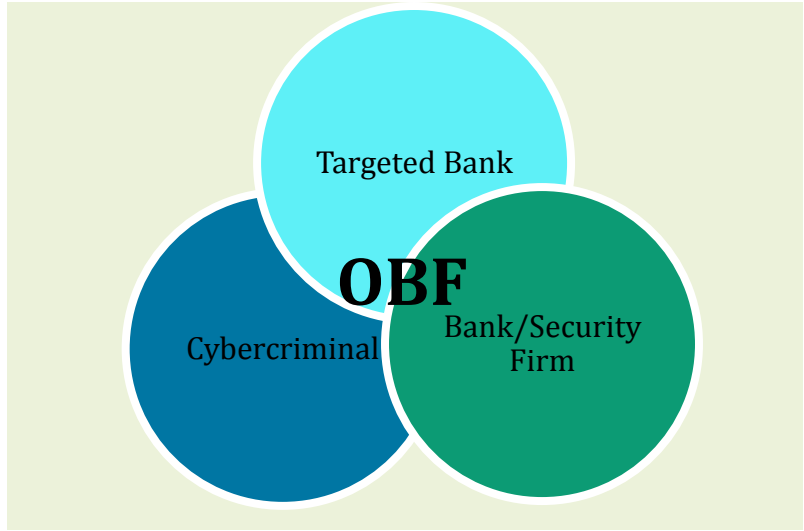


FIGURE 8- ONLINE BANKING FRAUD USING RAT FRAMEWORK

2.8 BUILDING A CONCEPTUAL FRAMEWORK

2.8.1 A RESEARCH GAP

Up to now, the origin of online banking and today's problems that are encountered by banks, security firms and users were delineated. We explained that as banks developed their online banking platforms to gain cost-efficiencies, the unique characteristics of Internet led to the raise of online banking fraud. Also, we argued that developments of security products such as antivirus software and security measures provided by security firms only mitigate online crime to a certain extent, and the problem is not only technical, but also socio-economic deficiencies and misaligned incentives of actors involved in the field augments the problem of online banking fraud.

In this research, we wish to address online banking fraud. But the question is which part of the online banking value chain should be address? In light of the three categories of Routine Activity Theory, most of the existing research and investments have focused on the 'absence of capable guardian' category. The reason is because this is the only part that is completely controlled by defenders, is known to them, and associated developments in this field include producing software products and services that can be sold for finical gain.

The knowledge gap exists when it comes to the other two criteria of RAT. The fact is that investigating the other two criteria requires in-depth knowledge about records of previous attacks, behavior of attackers, specifications of the attackers, and contains a certain level of uncertainty because it is not determined by defenders but by the attackers. Due to the typical problems researchers normally bump into for conducting research in the security field, such as lack of aggregated and accurate data, lack of access to attack datasets due to confidentiality reasons, lack of essential expertise relevant to the security field etc., these two categories either have not been studied much, or not fully analyzed.

The example of lack of reliable information about 'attractive target' category of RAT can be seen, considering a number of security reports such as threat report published by F-Secure (2012) in which a list of the top attacked targets in 2012 in online banking attacks via the SpyEye malware is provided. The reported list is not a completely valid representation of top attacked targets in 2012 because the ranks attributed to the attacked domains are calculated by simply counting the number of times that the domain name is seen in the malware specific files that contain information about attacked targets¹³. However, the number of times a domain is seen in the configuration files is not identical to the number of times it has been attacked but may happens according to several other reasons. In short, there exists no reliable information about 'target selection' by cybercriminals whereas some information in regards to the attacked targets may be found Malware configuration files.

2.8.2 THE RESEARCH'S UNIQUE OPPORTUNITY

As it is explained above, finding and accessing reliable data is one of the most important issues that researches often face in the field of cybersecurity. However, in this research, a unique opportunity of accessing a Zeus malware configuration files exists. Zeus is one of the most important widespread malware that is first exposed on July 2007. It mostly works on computers using Microsoft Windows operating systems and, since 2012, on Blackberry and Android phones. The Zeus configuration file contains all of the command and control scripts to tell the bot how to connect to the botnet and how to execute the attack on the mentioned domains (targets)¹⁴.

¹³ Malware configuration files. The specific characteristics of a configuration file will be explained in more details in the upcoming chapters.

¹⁴ More details about Zeus malware will be discussed are coming up chapters.

The Zeus dataset is provided by Fox-IT for Delft university of Technology. Fox-IT is a security company located in the Netherlands, delft that provides security solutions for more than 20 countries worldwide.

The opportunity of having the dataset provided by Fox-IT is unique primarily due to the reasons that are explained above. That is, because normally finding information about a malware configuration file which may contain confidential information about different companies and financial institutions is hard for independent researchers. Moreover, the dataset covers Zeus records of wide variety of countries and the records are gathered from 2009 until first quarter of 2013.

2.8.3 FORMULATING THE RESEARCH QUESTION

Among the two criteria of RAT that are studied less than others as its discussed in the previous section, '**Suitable Target**' is determined as the more plausible topic as for this research due to the following reasons:

- A research gap exists when it comes to reliable information about 'target selection'. A number of security reports have already performed similar research to determine highly attacked domains by analyzing a set of malware attack records. The already interpretations from the raw numbers in the dataset may not be necessarily correct which leaves an opportunity for more research in this field.
- The above gap nicely fits with a unique opportunity that exists for this research. The required scientific contributions for performing research on 'target selection', i.e. records of **Zeus** financial malware configuration files and expertise are provided by Fox-IT for Delft University of Technology for the purpose of scientific research. Although the dataset provided would needs to be aggregated for further analysis, it still counts as a rare opportunity and as a valuable input for an empirical research.
- Recent reports provided by security firms like TrendMicro and McAfee indicate that instead of broad generic attacks, financial malware attacks have moved towards being more target-specific, as can be determined from the recent online banking incidents (Sherstobitoff, 2013; TrendMicro, 2013). So, it is probable that by investigating financial malware records, a line of intelligence in selection of targets by cybercriminals can be found.
- Considering the fact that online attacks are executed by financial malware through instructions and targets provided by attackers that are in configuration files of malware, it is probable that some intelligence in regards to target selection could be found from attackers instructions and selected targets mentioned in financial malware configuration files.

This leads us to the following research question:

Research Question: *Can we extract intelligence on criminal attack patterns and target selection from the files which financial malware use as instructions for their operations?*

To answer this research question, we would need to answer a set of sub questions:

SubQ1: *a. What are the current most important online banking threats?
b. What is Zeus malware and what are its most important variants? How do they work in online financial attacks?*

SubQ2: *What are the problems of the Zeus malware dataset? Where are the gaps in the datasets?*

SubQ3: *Is it possible to build a dependent variable to explain patterns of target selection in online banking attacks discerned from the instructions available in the Zeus malware data?*

SubQ4: Is it possible to identify some explanatory factors that can explain the attack patterns and target selection by cybercriminals?

Please note that sub question 1a has already been answered as the part of the literature review; the rest of the sub questions will be answered in the upcoming chapters.

2.8.4 THE CONCEPTUAL FRAMEWORK

During this chapter we emphasized several times that that because cybersecurity is a relatively new problem, there is not always a clear consensus about norms and factors involved in the problems and their causal relations. In the current literature, something like a 'firm' hypothesis, i.e., clearly stated hypothesis that could be simply copied into a conceptual framework could not be found.

However, it is possible to generate an initial sketch from the relations based on the literature that is reviewed in this chapter and the factors that influence online banking fraud and more specifically target selection. The below schema (Figure 9) is the primary conceptual framework based on RAT model. This model is only the primary framework for determining the overall structure of the research. Figure 9 is further explained as below:

Our perception of 'Target'

Primarily it should be mentioned that although online banking malware attacks such as those by Zeus are executed ultimately through users' computers, attackers do not target specific user or users because of their specifications. Indeed, we believe that **'banks'** are the actual target of attacks due to their specifications, while **'users'** and their demographics are counted as properties of target banks.

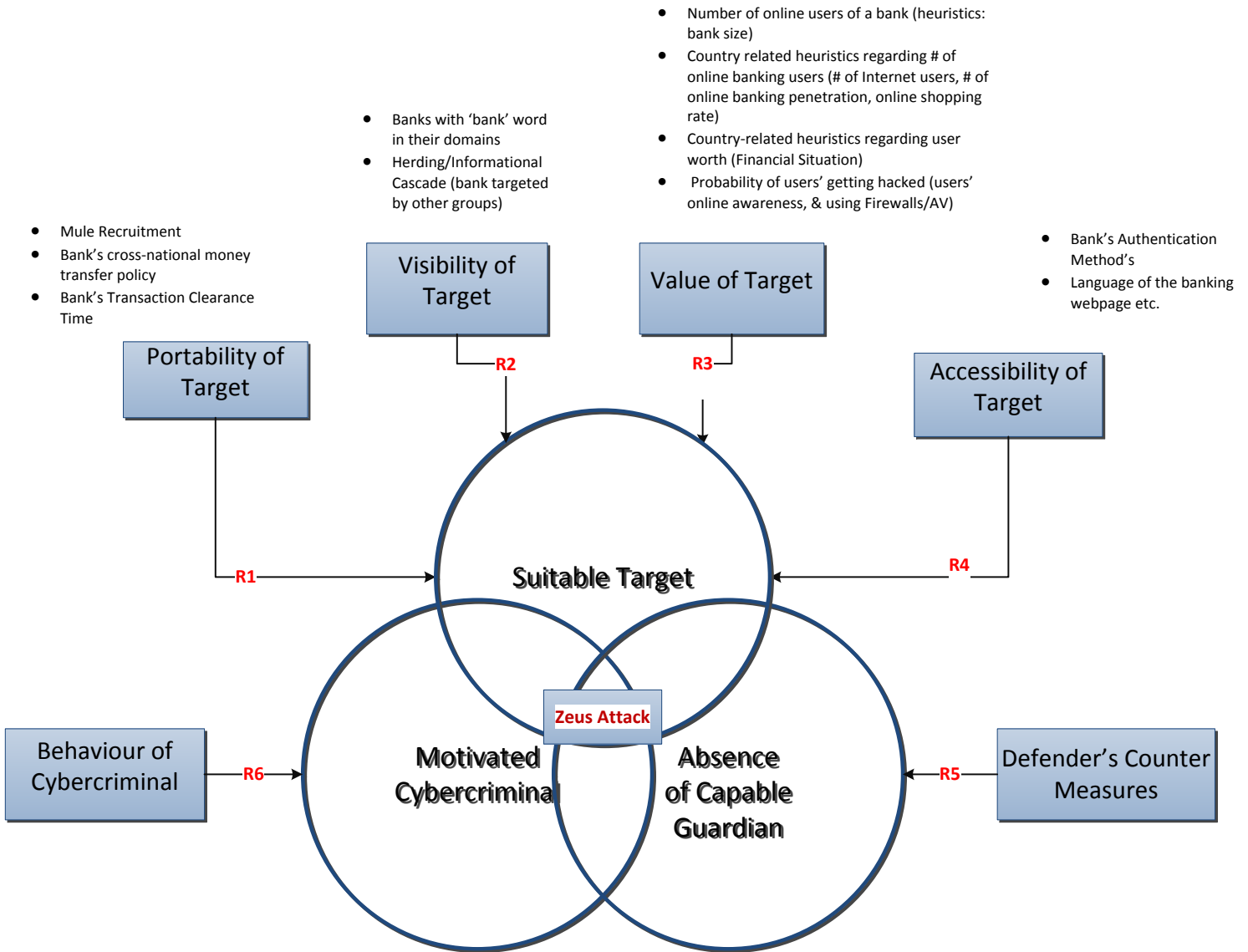


FIGURE 9- INITIAL CONCEPTUAL FRAMEWORK

Figure 9 displays the initial conceptual framework, built upon the points discussed in section 2.7.2. Referring back to the RAT approach, we discussed that crime occurs when the three elements namely suitable target, motivated offender and absence of capable guardian meet at the same time and in the same space (Cohen & Felson, 1979). In the framework, the three criteria are derived from the RAT approach. Likewise, an online banking attack by Zeus malware can occur in the common area of the three criteria.

Table 4 lists all of the major relations present in the conceptual framework. We also include a number of examples for each of the influencing factors in the table. Starting from 'suitable target', as stated by Yar (2005), suitability of a target for attack can be estimated due to its four-fold essential properties: portability (inertia), visibility, value and accessibility. Thus, a cybercriminal selects her target based on the above for properties.

Although the focus of this research is on the target selection concept, the behavior of a motivated offender and the absence of capable guardians are not negligible. It should be noted that in examples introduced in the 'Absence of Capable Guardian' category it is assumed that defender's counter measures are not known to

cybercriminals. It is mainly because information about banking security measures is normally confidential and not public and thus it is likely that attackers do not have access to none of such information. Therefore, it does not directly influence the decision of cybercriminal for selecting a target but ultimately can influence the occurrence of the crime itself. That is why defender's behavior, activities and counter measures are categorized in the 'absence of capable guardian' group.

For instance, If an attacker knows that his target has security by Fox-IT, the accessibility and thus the suitability of target from the viewpoint of the attacker is decreased whereas if he knows that Fox-IT is doing the security and he is also able to perform his attack as a criminal, then that is a control measure by the bank, and there is a chance that cybercriminal get caught (actions being done during or after the crime). Suitability of target is being varied if criminal knows about security before the crime.

Criteria	Main Factors Influencing a Target to be selected	Example
Suitable Target	R1: Portability of Target (Weight and size of the target that affects how fast the crime can be executed)	<ul style="list-style-type: none"> • Mule recruitment • A bank's cross-national money transfer policy • A bank's transaction clearance time
	R2: Visibility of Target (How visible the target is to the offender)	<ul style="list-style-type: none"> • Banks with the word 'bank' in their domains • Herding/Informational Cascade (bank targeted by other groups)
	R3: Value of Target (The value that can be gained by the offender if the attack occurs, determined by: <ul style="list-style-type: none"> ○ The number of online users ○ The amount of money that can be gained from each user 	<ul style="list-style-type: none"> • Number of online users of a bank (heuristic for bank size) • Country related heuristics regarding # of online banking users (# of Internet users, # of online banking penetration, online shopping rate) • Country-related heuristics regarding user's worth.
	R4: Accessibility of Target (How easy a target is to be reached)	<ul style="list-style-type: none"> • A bank's authentication mechanism • The language of the banking webpage
Absence of Capable Guardian	R5: Absence of Defender's Counter Measures (Unknown Measures)	<ul style="list-style-type: none"> • A bank's security measures • 3rd party real-time security measures for banks (Fox-IT DetACT) • Country related statistics about the probability of users' getting hacked (users' online awareness, & using Firewalls/AV)
Motivated Offender	R6: Behavior of Cybercriminal	<ul style="list-style-type: none"> • Cybercriminal's rational decisions based on target specifications and cost/benefit analysis

TABLE 4- LIST OF FACTORS THAT INFLUENCE ONLINE BANKING CRIME

Routine activity approaches suppose that actors are free to choose their courses of action, and do so based on anticipatory calculation of the utility or rewards they can expect to flow from the chosen course. Likewise, as it is noted in the Table 4 based on the opinion of security experts, behavior of 'Motivated Offender' or what we call it cybercriminal in online banking crime, is to large extent rational and follows 'Rational Choice Theory' introduced by Clarke and Cornish (1987). Rational Choice Theory argues that 'an individual acts as if balancing costs against benefits to arrive at action that maximizes personal advantage'.

However, we should consider that in the case of online banking fraud, no perfect information is available for cybercriminals in regards to their targets and the counter-measures that are being taken by offenders. So we guess that a cybercriminal's decisions follow Simon (1991) bounded rationality and mainly fits within one of the two categories below: They either use some heuristics (e.g., based on the target's properties) to select the target or they copycat each other's targets based on the previous successfully executed attacks. The latter is what we called informational cascade as a form of herding behavior previously discussed in this chapter. A point of caution should be raised that, both of the decisions above are considered partially rational because they are made based on the associated cost and benefits assigned to the attack, but only the information about both of them are bounded due to the lack of information and consensus about norms in the cyber security field in general.

Chapter 3 - Methodology

Introduction

In the previous chapter, we presented the research problem, the research questions, and the conceptual framework that guide us through the rest of this thesis.

In this chapter however, we are going to explain how the main research question and its subsequent sub questions will be answered in this research. In order to define the methodology that is going to be used, it is essential to recall the main research question:

Research question: *Can we extract intelligence on criminal attack patterns and target selection from the files which financial malware use as instructions for its operations?*

As the research question above implies, we are going to perform is a **quantitative empirical** research with an **explorative** objective (Kothari, 2009). The research is quantitative because it based on measuring quantitative malware dataset and it is empirical because it is a data-driven research and any conclusions in the research can be verified by experiments and observations.

The purpose of the research is explorative because we plan to explore the available empirical malware data to see whether it is possible to gain any insight about target selection by cybercriminals in online banking fraud by cybercriminals. As Shields and Tajalli (2006) explained, exploratory research is loosely coupled and mostly link to a conceptual framework built upon a number of expectations which could be the basic guide for preliminary investigations. In our research, the preliminary conceptual framework is built in chapter two. This framework will be our guide for doing the core of the research in chapter five, which is extracting a method in terms of defining a set of variables for explaining target selection by cybercriminals from malware data. Later in chapter six, among the variables defined in chapter five the one that matches our research question the most will be statistically tested along with a number of independent variables as for the proof of the concept introduced in chapter five.

Figure 10 provides an overview of the steps that are going to be taken in this research. Each individual step will be explained in detail in the upcoming sections.

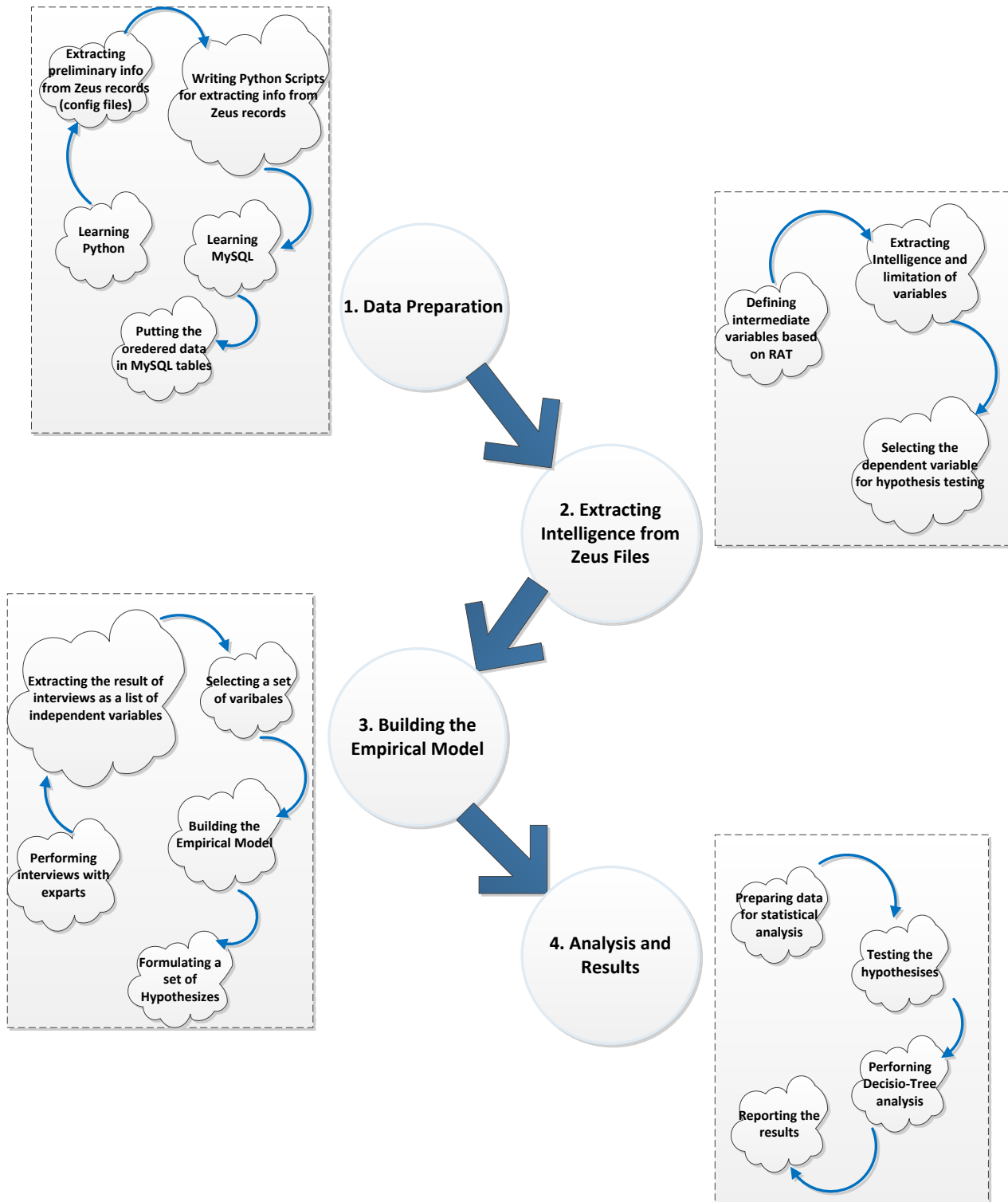


FIGURE 10- OVERVIEW OF THE RESEARCH STEPS

3.1 DATA PREPARATION

The data that is going to be used in this research is provided by Fox-IT, a reliable security company that works with trusted partners in more than 20 countries, headquartered in Delft, the Netherlands. The data consists of 11,000 samples of Zeus malware records of different online banking incidents in a time span of four years (2009-2013). The details of what Zeus malware is and how it works will be explained in chapter four of this report.

The Zeus records are raw and contain a lot of irrelevant information to our research. Thus, the first step of the research would be to convert the raw data into analyzable collection. In short, the Zeus malware records need to be processed, treated and aggregated. This task requires a lot of programming using **Python** scripts. Considering the number of records, with different formats, limitations and mistakes, and the fact that programming is not specialty of most of the students in TBM faculty, this part of research required great effort time-wise.

Data preparation in itself has different steps. First of all a number of samples from the configuration files will be explored to find out what possible information can be gained from these files. Next, according to our research question, a preliminary selection of relevant information will be made among all of the info available in the Zeus records using **Python** scripts. In the last stage of the data preparation, extracted information from the configuration files will be stored in MySQL tables for easier access in later stages.

3.2 EXTRACTING INTELLIGENCE FROM ZEUS MALWARE CONFIGURATION FILES

This step is the core of the thesis and the main part of the research question will be answered in this chapter. After placing all the required preliminary information from Zeus malware records into set of MySQL tables through the data preparation step, the explorative phase actually starts.

In this phase, according to the conceptual framework introduced in chapter two, a method for extracting intelligence from Zeus malware data will be defined in form of a list of intermediate variables. These variables will be explored one-by one in terms of the extent that they can explain target popularity among cybercriminals from Zeus records and the possible intelligence and limitation that they can provide. At the end of this chapter the variable that is able to address our research question the best; that is, the variable that would be able to explain target popularity among cybercriminals the best would be selected to be used in later steps for proving the method that is introduced in chapter this chapter.

3.3 BUILDING THE EMPIRICAL MODEL

As it is explained above, in chapter five we have already answered the most important part of our main research question by suggesting a method for extracting intelligence from Zeus data. Next we intend to shortly examine the dependent variable introduced in chapter 5.

For this, firstly, we need to come up with a set of independent variables that may influence the dependent variable. These variables should be found from either literature or opinion of experienced people in our case security experts. Due to existence of very few information in the literature about target selection by cybercriminals in online banking fraud as it is discussed also in the introduction chapter, we gathered our independent variables by the means of conducting four set of semi-structured interviews with three security experts from Fox-IT and one expert from criminology field. Among all of the independent variables identified those that empirical data for which the empirical data will be accessible. It should be noted that because this section of the research will not

going to be one of the main contribution of the author and is conducted only as a proof of concept, it will not contain a comprehensive set of all possible independent variables but only a few.

After identifying the variables, an empirical or empirical model will be built and the relations in the model will be converted into a set of five, empirical hypothesizes for further analysis in the next steps.

3.4 DATA ANALYSIS

The analysis of the relations would remain as the last step of this research. This step will be done through employing a number of statistical analysis tests such as bivariate analysis and comparison of means for testing the relation between dependent and independent variables using **SPSS** software. Rather than empirical hypothesis testing, this section will also contain a number descriptive statistics about other intermediate variables that are identified in chapter five. These descriptive statistics will provide information about target selection from different perspectives.

The dependent variable will be analyzed using a data mining method as well. A very brief decision-tree classification test will be performed to understand whether any patterns could be found in the target selection by cybercriminals from the Zeus malware data based on number of attacks that are executed on each domain. The dataset that is going to be used in the test will contain the number of attacks on each domain as the decision variable and three attributes as , continent of the attacked domain, domains with English option in their banking webpage and domains with word 'bank' in their domain names. The decision tree test will be performed using **WEKA** software.

Chapter 4 - Data Preparation

Introduction

Recalling from chapter three, the first step of the research is to aggregate the raw dataset that is provided for us. In short, the data can only be used for exploration and analysis steps if it has a certain level of aggregation. In our case, taking into account the volume of information available in each malware configuration file that may not necessarily be relevant to our work, the raw data should be parsed, and the relevant required information should be extracted and stored in a structured database such as MySQL tables. Of course, the relevancy of information will be determined by the scope of the research.

However, before going into more details, first we have to know what our dataset is and what specifications it has and what limitations it would possibly impose on our research. To be able to address the above notions, we have to answer the following sub questions in this chapter:

Sub Question 1: *What is Zeus malware and what are its most important variants? How do they work in online financial attacks?*

Sub Question 2: *What are the problems of the Zeus malware dataset? Where are the gaps in the datasets?*

It should be noted that the sub question 1a is already answered in chapter two and the above two sub questions will be answered in this chapter. In the following sections, we will address the type and origin of the dataset that is going to be used in this research. Considering the type of our dataset that contains a number of Zeus-variants configuration files, a comprehensive explanation will be provided about the similarities and differences among all the variants of the Zeus malware. Further, the process of aggregating the data as well as extracting the required primary information from the data will be explained and discussed.

The research of this thesis is based on a dataset that contains configuration files of Zeus malware and its variants provided by Fox-IT for Delft University of Technology. To understand the datasets and to be able to aggregate them first we need to know what Zeus malware is and how does it work.

4.1 ZEUS BANKING TROJAN MALWARE

Zeus also known as Zbot is a readily available and the most widely-spread malware package contains required tools to build and control a botnet. Zeus first exposed on July 2007, mostly works on computers using Microsoft Windows operating systems and, since 2012, on Blackberry and Android phones. It is being sold and traded to many customers in underground forums while the kit is very simple to use since it does not require too much technical knowledge (Wyke, 2011). This resulted in a huge number of independent Zeus-buyers who create their own botnet which indicates that the numbers of Zeus botnets are quite large. The Zeus malware is primarily designed for stealing online banking information but according to its features it also has been used in other types of data or identity theft (e.g. password sniffing).

4.1.1 ZEUS TECHNICAL FACTS

The Zeus kit contains:

- A control Panel application that is used to maintain/update the botnet, and to retrieve/organize recovered information.
- An EXE builder that is used to create the Trojan binaries and encrypt the configuration file.

Each buyer/customer of the Zeus kit has to build his own bot executable, the ones that he intends to distribute to his victims. The first step in building a bot executable is to edit the configuration file. Thus, the builder is used by

each customer to create both the encrypted configuration file and the bot executable that is specific to the customer in terms of location or the configuration file and RC4 key (password) (Macdonald, 2011).

Zeus host contains three components:

- a configuration file (being encrypted with file extension *.bin)
- an executable file that contains the newest version of the Zeus Trojan
- a drop server (mostly a PHP file)

CONFIGURATION FILE

The configurations file contains all of the command and control scripts to tell the bot how to connect to the botnet. Moreover, it contains information on which users to attack, what user data to gather and how to do so. The configuration file has two parts:

1. Static Configuration

There is a builder tools available in the Zeus kit and the StaticConfig could be compiled into the binary by the builder tool. It contains information that the bot will need when it is first executed.

The file mainly contains:

- The name of the botnet that this bot belongs to.
- The URL where the bot can get the dynamic configuration file.
- The encryption key that is used to hide information transmitted within the botnet (password).

2. Dynamic Configuration

The DynamicConfig is downloaded by the bot immediately after it is installed on a victim's computer. This file is downloaded at timed intervals by the bot, and can be used to change the behavior of the botnet. Most of the entries control how information is collected from the infected computer (Macdonald, 2011).

Available settings include:

- A URL where the bot can download a new version of itself, if the command to do so is given.
- The URL of the drop server where logs, statistics and files will be uploaded and stored.
- Information used to inject additional fields into web pages viewed from the infected computer.
- A list of URLs where an emergency backup configuration file can be found.
- A set of URL masks used to cause or prevent logging of information.
- A set of URL masks to indicate that a screen image should be saved if the left mouse button is clicked.
- A list of pairs of URLs that are used to cause redirection from the first URL to the second.

Data sent through the Zeus botnet is encrypted with RC4 encryption. In this implementation a key stream is generated from the botnet password, and is XORed with the data. To encrypt the data that is communicated in the botnet, the same password is used (Macdonald, 2011).

EXECUTABLE FILE

The EXE file built by the builder component is to be deployed by the botmaster. The EXE file is unique for each customer, in the sense that the location of the configuration file embedded into the binary and the RC4 password is different even if they use the same version of Zeus.

So the only thing that differentiates one Zeus kit created botnet with another is the configuration details while the functionality and the behavior will always be the same.

DROP SERVER (DROP ZONE)

This is the server component of the Zeus kit and has a collection of PHP scripts that allow the botmaster to monitor the status of their bots, issue commands to them and retrieve the information that are collected by bots. The data stolen by the bots is also sent to the drop server.

HOW DOES ZEUS OPERATE?

The first step of the Zeus malware infection is that the victim's computer gets infected. In most of the cases this happens by clicking on spam emails containing links or infected files (Macdonald, 2011). The infection flow of the Zeus botnet is illustrated in Figure 11.

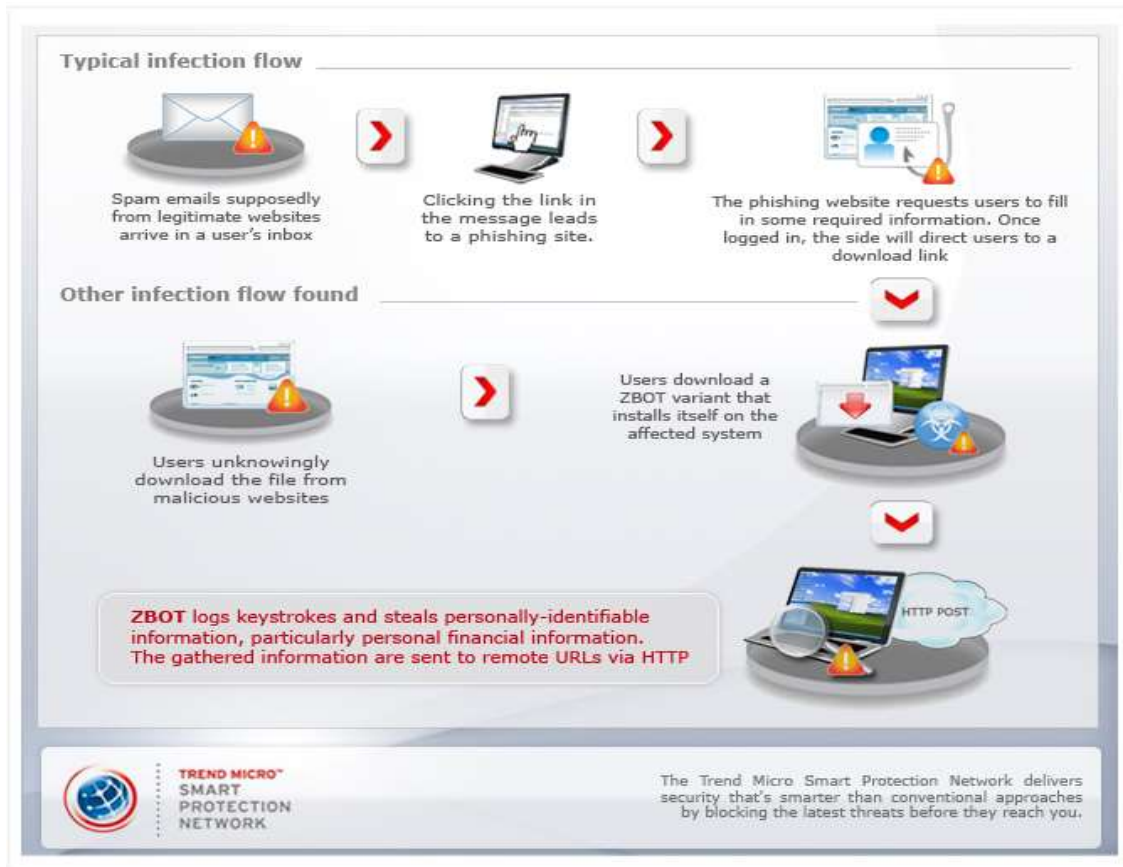


FIGURE 11- ZEUS INFECTION CYCLE (TRENDMICRO, 2010)

When a computer gets infected by the Zeus Trojan malware and become a bot, the bot communicates with the C&C server and requests a dynamic configuration file. Immediately, dynamic configuration files being encrypted by the C&C server is being send to the bot. As soon as the configuration file is received by the bot, the bot will retrieve the drop server URL from it and will HTTP POSTs information about itself to the drop server and updates its status to 'online'. It monitors the user's Web browsing activities (both HTTP and HTTPS) using the browser window titles or address bar URLs as triggers for its attack and sends the logs on average every minute to the server(Wyke, 2011).

The server then, starts to communicate with its bot to send the commands or inject JavaScript codes into web pages of the infected computer. This is being done on-the-fly, as data is passing from the server to the browser of

the client. Other ZBOT variants display a second fake login page after the original login page to get additional information. The following script is an example of injection to a user's browser:

```
set_url http://www.bank.com/login.html GP
data_before
name='password' *</tr>
data_end
data_inject
<tr><td>PIN:</td><td><input type='text' name='pinnumber' id='pinnumber'
/></td></tr>
data_end
data_after
data_end
```

Targeted URL

Injected Scripts

The *set_url* parameter is the page to be attacked; *data_before* contains the text to search for before the injection point and *data_inject* has the text that will be injected. Figure 12 is an example of a login page before and after injection.

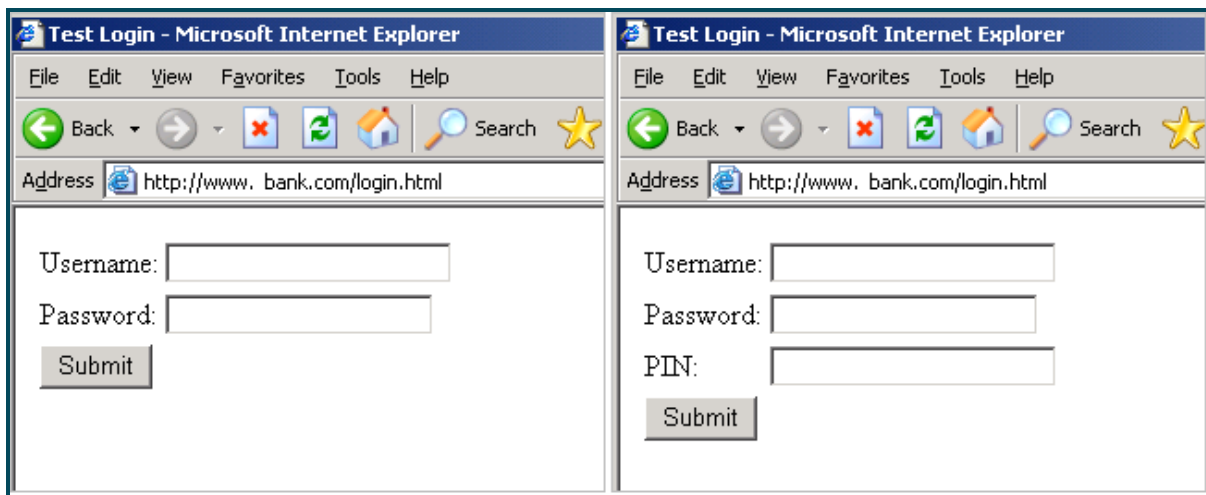


FIGURE 12- EXAMPLE OF A WEBPAGE AFTER ZEUS WEB INJECTION (MACDONALD, 2011)

As it can be seen from the above example, Figure 12, an extra field is injected into the original banking form. As soon as the victim enters his information (banking credentials, credit card number etc.), the data is captured by attackers.

Zeus control panel has a handy feature, called 'Enable No-Shit reports,' which only select predetermined specific information from all of the stolen data by the malware such as 16-digit credit card numbers and data that victims are submitting to pre-selected online banking sites (Krebs, 2011).

WHY ARE ZEUS THREATS PERSISTENT?

In addition to its social engineering tactics and **EVER**-evolving spamming techniques that are used by Zeus malware attacks, ZBOT makes detection difficult because of its rootkit capabilities. Upon installing itself on an infected system, Zeus creates a hidden folder in the system directory to prevent itself from being discovered and removed by users. Also, Zeus is able to deactivate Windows Firewall or dismisses itself in case a particular firewall processes start working on the system (Macdonald, 2011).

Also, Zeus botnet updates its feature with trends in operating systems. Newer variants provide full and integrated support for other Windows operating systems like Vista and Windows 7.

WHAT IS THE DIFFERENCE BETWEEN VERSION 1 AND 2?

The main improvement of the second version Zeus compare to the first version is its support for the obfuscation of the RC4 key and an extra level of encryption on top of RC4. In the second version also the configuration data is stored in the registry rather than a file on the disk. The easiest way to retrieve the RC4 key and URL to download the configuration file is from the Zbot PE file that is injected into running processes. Once the configuration file has been RC4 decrypted there is an extra XOR decryption on top.

In May of 2011 the source code for version 2.0.8.9 was leaked onto the Internet, resulted in creation of new versions, or derivatives of Zeus malware by other authors. The mutual point of all derivatives is that they all are built based on the Zeus source code. Up to now, three major Zeus-family malwares are:

- Citadel
- Ice IX
- Peer-to-peer (P2P) version

As reported by TrendMicro, Figure 13 displays the distribution of the original Zeus and its variants during April-May 2012. As it can be seen from Figure 13, Citadel has the majority of the market share in 2012 compared to other variant of Zeus.

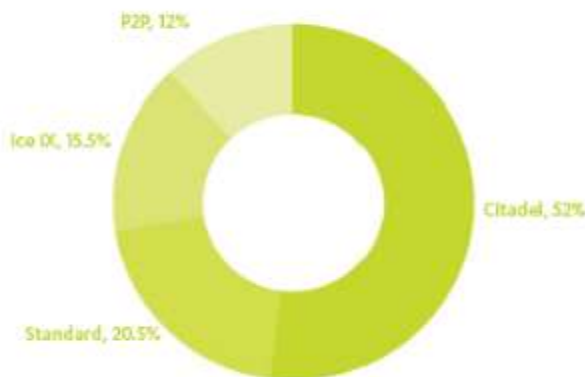


FIGURE 13- ZEUS VARIANTS' DISTRIBUTION (TRENDMICRO, 2010)

4.1.2 CITADEL

Citadel is one of the variant of Zeus Trojan that has become active from early 2012. Citadel has other usages than just stealing banking credentials since a number of targeted attacks, especially of government organizations have been observed in 2012 using this malware. This usage goes beyond the norm of financially motivated crime-ware (Sherstobitoff, 2013). According to a McAfee report about Citadel, the targets of this malware are chosen selectively based on specific patterns. The victims are primarily in Europe; Figure 14 displays the distribution of the Citadel malware across Europe. As it is obvious from the picture, United Kingdom, Germany and the Netherlands have been countries of great interest to cybercriminals. Figure 15 displays the citadel distribution worldwide. The red circles displayed in Figure 15 are indicators of the locations that have been attacked by Citadel from Dec 2012 to Jan 2013.

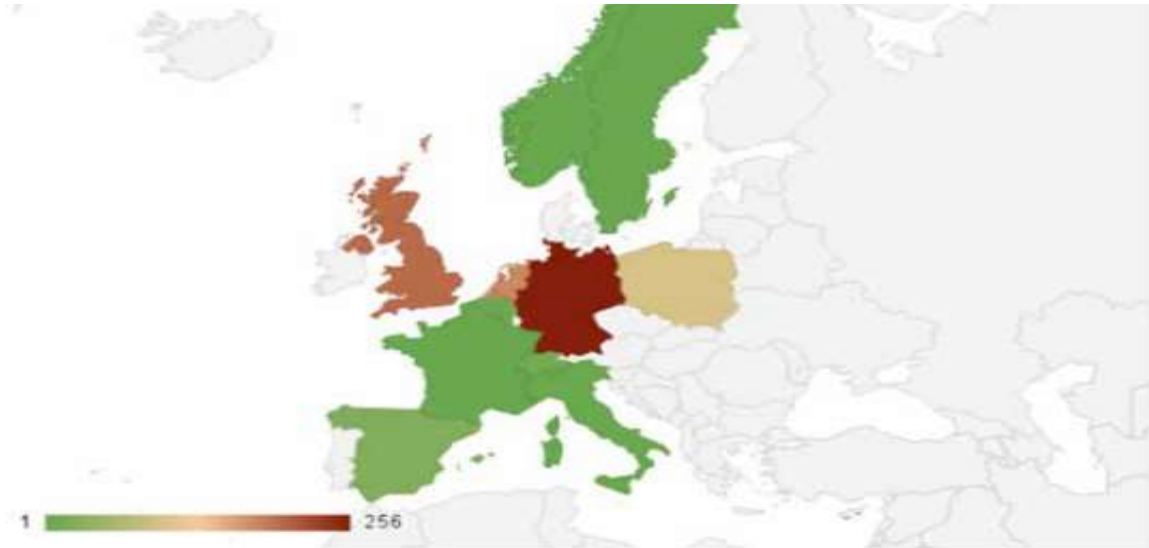


FIGURE 14- CITADEL DISTRIBUTION IN EUROPE (SHERSTOBIKOFF, 2013)



FIGURE 15- WORLDWIDE CITADEL DISTRIBUTION (SHERSTOBIKOFF, 2013)

Citadel Trojan is said to be more organized than Zeus in terms of executing the online attack and has features that help it to extend beyond targeting customers of financial institutions such as features that allow it to collect anything from a victim's PC (Sherstobitoff, 2013). Ahnlab's Malware Analysis report published on December 2012 mentioned that Citadel matches original Zeus by approximately 75% while the remaining 25% consists of the new features that are added uniquely in Citadel (AhnLab, 2012). The new features mainly are (RSA, 2012):

- **Local Pharming:** This feature allows the C&C server to decide which URL can or cannot be reached by the victim and which page should the victim lands if he/she tried to reach the blocked URLs. This feature is placed in configuration files in a section called 'DNS-filters'. Any attempts for reaching URLs listed there would be redirected to a unique IP address (e.g. <http://74.125.224.72/> for google.com). This feature enables the cybercriminal to both create more 'reliable' phishing attacks and isolate victim machines from Anti-Virus or police services.
- **More function hooks:** It covers a much larger range of Windows functions than Zeus did.
- **The C&C server side:** The Citadel botnet uses Zeus control panel but strengthened against web-based attacks that plagued Zeus. The panel has also been improved in terms of user interface and ease of use.

- **Trojan's encryption method:** Citadel added an extra layer of protection for communication between its bots and C&C servers, including Login key AES key.
- **Video-grabber:** The Citadel malware has video-grabber in the dynamic configurations files which lets the malware record the video of the victim during the attack.

4.1.3 ICE IX

Ice IX is one of the variants of Zeus developed based on the Zeus v2 source code. The Ice IX Trojan possesses improved capabilities of Zeus as well as several additional features that did not exist in the original Zeus. Apparently, the most valuable feature of Ice IX is the implementation of a defense mechanism designed to evade tracker sites such as Zeus Tracker. The dynamic configuration file cannot be downloaded directly but through the proxy.php file which can only be opened using the same key that is used for decrypting the static configuration file. If the request for configuration is created not by bot with the same key the 404 error will be returned. So there is no way to download and analyze the configuration file (RSALabs, 2012).

4.1.4 PEER-TO-PEER (P2P) ZEUS

P2P Zeus is a significantly improved version of the Zeus Trojan, because of its completely different communication structure. In P2P Zeus, the centralized C&C server, a single point of failure targeted by researchers and law enforcement is replaced with a robust P2P network (Symantec, 2012). In the P2P model of Zeus, each infected machine (bot) upholds a list of other infected machines. In this way all of the peers act an enormous proxy network between the P2P Zeus botnet operators and the bots. The distribution of configuration files, the propagation of binary updates and sending stolen data to the controllers are being done through the peers in the network instead of relying on the centralized, vulnerable C&C server (Stone-Gross, 2012).

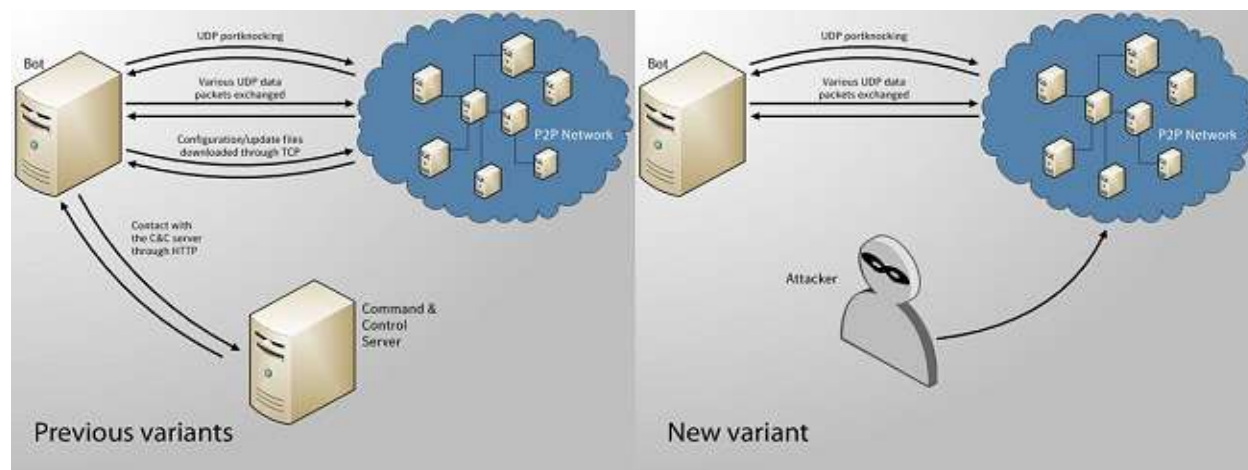


Figure 16- The P2P C&C paradigm has shifted from central server to botnet (Symantec, 2012)

4.2 OVERVIEW OF FOX-IT MALWARE DATASET

Fox-IT is one of the world's leading security companies located in the Netherlands. They are mainly active in preventing and mitigating online threats as a result of cyber-attacks, fraud and data breaches with solutions for government, defense, critical infrastructure, banking, and commercial enterprise clients worldwide. Accordingly, they not only have access to a vast quantity of information about infected systems of their customers worldwide, but also they gather real-time data of online fraud by manipulating systems as bots. The dataset has a set of unique specifications that are going to be explained in the below sections.

4.2.1 GENERAL SPECIFICATIONS

Table 5 presents the general specifications of the dataset used in this research. This dataset contains records of configuration files of Zeus and its variants financial malware mainly Traditional Zeus, Citadel, Ice IX and P2P Zeus. Table 6 displays number of domains that are attacked in each year based on our dataset. It should be noted that of course the number reported in Table 5 includes the commons domains in all four years.

TABLE 5- GENERAL SPECIFICATIONS OF THE RESEARCH DATASET

Terms	details
#of Configuration files	144625
#of text files	77899
#of text files contain target domains	10643
#of Botnets	2125
# Unique Domains	2225
Infection Date Range	21/1/2009- 4/3/2013

TABLE 6- NUMBER OF TARGETED DOMAINS IN DIFFERENT YEARS

Year	Unique domains
2009	1232
2010	1166
2011	1077
2012	1198

4.2.2 DATA COLLECTION METHOD

This dataset has been indeed collected by Fox-IT and its partners by running Zeus malware samples (executable binaries) files and creating a system that actually looks like a bot to emulating Zeus malware. By acting like a legitimate bot, the emulated system then is able to receive and download malware samples (binaries). Accordingly, Fox-IT retrieves and decrypts the malware configuration files by running the samples on their systems.

In order to receive updates of the configuration files which may contains new information for bots, this process is repeated regularly. As it is already explained in the previous chapter, bots in a Zeus botnet communicate with their bot-master through dynamic configuration files. The static configuration placed into the executable binary file, contains URLs from which the dynamic configuration file can be downloaded by the bot. In some of the cases, this dynamic configuration may also contain additional locations which contain additional dynamic configurations. In short, the process of downloading the dynamic configuration file is partially recursive. This feature is present to prevent a takedown of a single or several URLs or hosts leading to loss of control over the entire botnet by the bot-master.

4.2.3 DATA FORMAT

The format of the collected data can be recognized from Figure 17. The process of gathering data has been done in a structured manner. As soon as the emulated system downloads the encrypted binary dynamic configuration file, it starts the attempt to decrypt the files with all the available RC4 keys that have already been extracted by the system. In the cases that the key for decrypting the configuration file is found, it will be decrypted into a text file and its associated meta-data would be stored in a MySQL database called 'Zeus'. The dynamic configuration files (text files) contain all the actual commands and target domains (web-fakes and web injects) that was intended to be communicated to the bot from the C&C server.

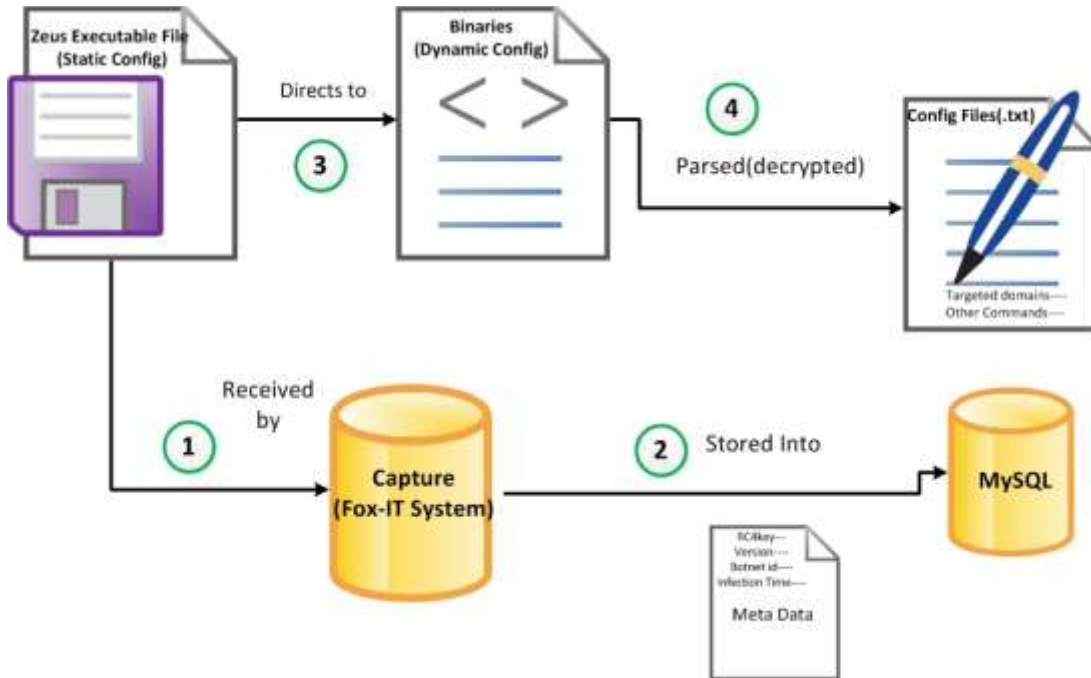


FIGURE 17- THE PROCESS OF AGGREGATING THE ZEUS MALWARE CONFIGURATION FILES

In order to provide more understanding of how the data is created, aggregated and used, the two set of files will be discussed below in more details:

1. Metadata (SQL Database)

Figure 18 displays the relation between the main tables in the MySQL Zeus database. The MySQL Zeus database is connected to the actual configuration text files via the name of text files. This means the contents of the 'md5' column in the 'files' table of the Zeus MySQL database is identical to the name of the configuration text files.

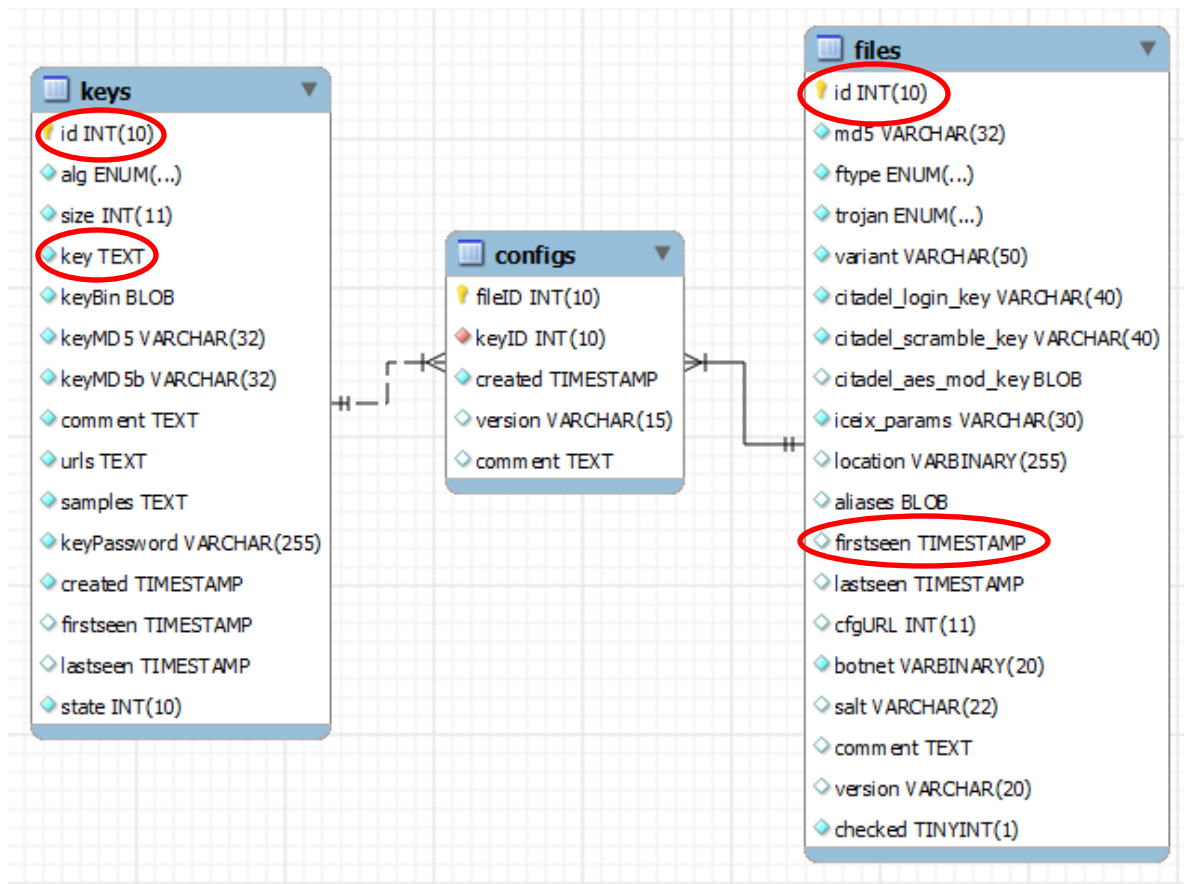


FIGURE 18- TABLES IN THE MYSQL ZEUS DATABASE (RED CIRCLES INDICATES THE VARIABLES THAT ARE USED FOR THIS RESEARCH FROM EACH TABLE)

As it can be seen from Figure 18, the 'Keys' and 'files' tables are connected to each other via 'configurations' table in the Zeus database. The first box displays 'files' table that mainly contains general info about each of the configuration files. The 'id', 'md5' (filename) and 'first seen' (Infection Time) columns of this table are mostly used for the purpose of this research.

The second box in the above figure is the 'configurations' table which acts as the connection between 'keys' and 'files' tables. The right table in the above figure shows the 'keys' table with its associated columns. From this table only id and key (key-state or RC4 key) columns are used for the purpose of this research. In short, **botnet id** and **infection time** are the most useful information that we extract from these tables for our research.

2. Configuration Text Files

As mentioned in the previous section, a configuration file is a file with which a botnet's command and control server communicates its commands to all of the bots in the botnet. These commands contain actions that should be taken by individual bots in a botnet including the scripts that should be injected into the user's browser, targeted domains (banks URLs) and other version-specific features like recording video and/or key loggers. In case of Citadel Trojan, due to its sophisticated features, the configuration files are dynamically updated and so commands from C&C server are communicated in real-time. Figure 19 displays how these configuration files look like.

```

WebInjects:

set_url */my.ebay.com/*CurrentPage=MyeBayPersonalInfo* <FLAG_GET><FLAG_LOG>
data_before
    Registered email address</td>*<img*>
data_after
    </td>
data_inject
    e-mail:

set_url *.ebay.com/*eBayISAPI.dll?* <FLAG_GET><FLAG_LOG>
data_before
    (<a href="http://feedback.ebay.com/ws/eBayISAPI.dll?ViewFeedback&*">
data_after
    </a>
data_inject
    Feedback:

set_url https://www.us.hsbc.com/* <FLAG_GET><FLAG_LOG>
data_before
    <table cellpadding="0" summary="page layout">
data_after
    </table>

```

FIGURE 19- PREVIEW OF ZEUS MALWARE CONFIGURATION FILE

4.3 BUILDING THE RESEARCH DATASET

In the previous section, an overview of the specifications of Zeus malware dataset is provided. In this section we are going to explain the process of extracting the main information from the whole set of parameters that could be found in the malware configuration files.

4.3.1 EXTRACTING THE KEY VALUES

Before starting to build the independent variable(s) and their accompanying metrics, we initially need to extract the information or values that we intend to use among all of the information available in the Zeus malware configuration files. The reason why we are extracting the values below is primarily to define our scope of work out of all the information exists in the data files and also to direct our analysis toward finding the most relevant dependent variable(s) to our work in this research.

TABLE 7- PRELIMINARY INFORMATION EXTRACTED FROM THE MALWARE CONFIGURATION FILES

Value	Details
Domain name	The targeted domains
Infection Time	Time of attack
RC4 Key (Botnet key)	Botnet's unique identifier

In the following section we will explain how the above values can be extracted from our dataset.

DOMAIN NAME

The format and content of the configuration text files are displayed in Figure 19. Each of the files contains the following values interesting for us:

- Name of the configuration file(text file)
- List of the targeted domains

The configuration data revealed that a total of 14870 unique URLs were targeted for web-inject. Not all of these URLs included the domain names. Many domains had several different URLs leading to them, using different paths. Sometimes, only the path is used for identifying a targeted website. Accordingly, extracting the domain names from URLs was not very smooth. As it can be seen in Table 8, most of the URLs that are included in the configuration files by cybercriminals are the obfuscated versions of the real URL. Clear from the below table, in several cases, cybercriminals use special characters like '*' to hide a part of parts of the URL. This means that, in several cases the domain name could not be found directly by having the URL. Other examples of such wildcards in the targeted URLs can be found in the below table. After excluding URLs with missing domain names and duplicate domains, a total of 243 unique domains were left.

The first column displays raw URLs in the configuration files and the second column displays our extracted domains associated with the URL in the first column. The question is how we extract the domain name from the wildcards in the targeted URLs.

TABLE 8-EXMAPLE OF TARGETED URLs IN THE CONFIGURATION FILES

Dynamic URL	Actual Domain Name
https://*/onlineserv/cm/index.cgi?state=otpsignin	afsbonline2.com
http*://*accountcentralonline.com/cmuser/login*	accountcentralonline.com
/logon/challenge/un/account*/card*	abnamro.nl
*https://*mybank.alliance-leicester.co.uk/view_accounts/*	alliance-leicester.co.uk
http*://www.*.*.banquepopulaire.fr*showportal.do*	banquepopulaire.fr

MECHANISM EMPLOYED FOR EXTRACTING TARGETED DOMAINS

In order to be able to extract the domain names from configuration files in the dataset, we used Python Scripts. However the procedure was not straight forward. In the below section the steps taken to extract domain names is from configuration text files are explained in more details.

I. Creating the URL text file

The goal of the first part of the procedure is to build a table that contains URLs and their associated filenames. The result of the below script has also put in a table in MySQL database. In order to be able to extract the above values from configuration text files, the algorithm presented in the figure below is performed:

1. Extract the name of files and put them in a list consists of filenames in memory and for each filename remove ".txt" extension from its end.
2. Loop through each line of each configuration text file and do the following:
 - i. Find the line contain "set_url" or "Target URL" keywords
 - ii. Extract the targeted URL from the line
 - iii. Determine the keyword the URL is extracted from e.g.: set_url or Target URL
 - iv. Connect the URL to its associated filename and put them in a list
3. Print the list as a text file with the below format
(Filename, URL, keyword_type)

BOX 1-GENERAL ALGORITHM FOR EXTRACTING URLS

The box below illustrates scripts written in Python for this purpose.

```
''' Created on May 15, 2013 @author: samaneh '''
import MySQLdb
import sys
import fnmatch
import os
import pprint
import heapq
import csv
from publicsuffix import PublicSuffixList

matches=[]
maindict = {}
psl = PublicSuffixList()

for root, dirnames, filenames in os.walk('c:/****'):
    for filename in fnmatch.filter(filenames, '*.txt'):
        matches.append([root, filename])

#get rid of filename's extentions
for root,filename in matches:
    filename_key = (os.path.join(filename).strip()).split('.',1)[0]

#counting the domain names seen per file and put them in a dictionary
list_url=[]

    fullfilename = os.path.join(root, filename)
    f= open(fullfilename, 'r')
    for line in f:

#if 'Target URL' in line:          # example: [Target URL           : 'https://www.barclays.es/*']
    if line.strip().startswith('Target URL'):
        url = line.strip().split(' ',2)[1]
        typ= 'T'
        list_url.append((url, typ))

#elif 'set_url' in line:          # example: [set_url https://onlineeast#.bankofamerica.com/cgi-
bin/ias/*/GotoWelcome]
    elif line.strip().startswith('set_url'):
        url = line.strip().split(' ',2)[1]
        typ= 'S'
        list_url.append((url, typ))

#combining the two dictionaries
    maindict[filename_key] = list_url
```

II. Extracting domain names from URLs

After all of the targeted URLs got extracted from the configuration files and their associated targeted domains were found, we stored them in a separate table in the 'Zeus' MySQL database. Figure 20 displays the final relation between newly generated and existing tables in the 'Zeus' MySQL database.

This step of the aggregating the data is one of the most time-consuming steps among all. In this step we intend to extract domains from the URLs we stored in our 'ds2_urls' MySQL table. As we explained before, extracting the domains from URLs is not straight forward due to the diversity and mistakes that are existed in the URLs' formats. Because the inputted URLs in the configuration files are not structured and they are free-format, every kind of formats can be existed. Moreover, because they codes are written by human-beings, many typing mistakes can also be seen in the URLs. In general text parsing is hard, and in this case it got harder due to the existence of obfuscated URLs and mistakes. The general outline of the Python script is provided in the box below.

The code follows three strategies to convert a URL to a domain

Strategy 1:

1. Tries to extract the domain part from URL:

If the URL be well-formed, then the domain can easily be extracted and the below three steps would be followed to extract the domain part from URL:

Step One: Removing URI

Step Two: Removing Sub-domain

Step Three: Removing the last part of URL (if any)

If the URL be obfuscated with wild-cards then the domain name will be extracted within strategy 2 and 3.

Strategy 2:

After performing the first strategy, a list of URLs will be made out of correctly extracted domains.

However, still a series of URLs could not be converted into domain names since they are too vague.

In this strategy, the remaining URLs from last section will be compared with our list using Regex matching, to see whether any match could be found for the URLs.

The matching occurs in two conditions:

If a match could be found, then the domain can be easily extracted because it only is matched with one URL.

If a URL get matched with more than one URL from the list, then the domain that matches the URL with matching probability more than 70% will be chosen. For instance, if a URL matches with 10 URLs and 7 of that contains "ebay.com" in them, then the code would match the URL with "ebay.com".

Strategy3:

30-40 of URLs are mapped manually because they are exceptions.

We are mapping them manually because including them in the scripts would make it complicated and hard to understand for the reader.

BOX 2-GENERAL ALGORITHM FOR EXTRACTING DOMAINS

The box below illustrates scripts written in Python for this purpose. As it can also be understand from the script, the code contains solutions for dealing with exceptions existed in URLs.

```

# Script to extract targeted domains from Zeus-Config-Urls
# 20-05-2013

import re
from publicsuffix import PublicSuffixList

simple_urls = {}
other_urls = {}
psl = PublicSuffixList()

def extract_domain(url):
# really stupid wildcards
    exceptions = {'exception URL': 'correct domain', ...}

    url = url.lower()
    if url in exceptions:
        return exceptions[url]

# STEP ONE: strip the https:// parts
    if url.startswith('http://') or url.startswith('https://') or url.startswith('http*://') or
url.startswith('https*://') or url.startswith('htt*://') or url.startswith('ht*://') or
url.startswith('http*//') or url.startswith('*http://') or url.startswith('*https://') or
url.startswith('*://') or url.startswith('*//'):

        # starting with ht and ending with double-slash
        # (the *// cases could be tricky, manually checked and are ok)

        assert '//' in url
        u2 = url[url.index('//')+2:]
    elif url.startswith('http*') or url.startswith('https*') or url.startswith('htt*') or
url.startswith('http://*') or url.startswith('https://*'):
        # starting with ht and ending with star
        assert '//' not in url and ':' not in url[6:]
        u2 = url[url.index('*'):]
    else:
        # remaining cases; should start with * (but those that don't are probably errors, but wont'break
anything)
        u2 = url[:]

# STEP TWO: strip some other wildcards from the start
    if u2.startswith('*'):
        # e.g.: https://*.westpac.com.au/esis/login/srvpage*
        u2 = u2[2:]
    elif u2.startswith('*') and u2[:2]!='*//':
        # Avoid stripping e.g. https://*/xxx ...
        # This can in some cases still cause problems, e.g.: *-sparkasse.de* (not handled yet)
        u2 = u2[1:]
    elif u2.startswith('www*.') or u2.startswith('www#*.'):
        u2 = u2[5:]

# STEP THREE: trim the end.
    u3 = None
    if len(u2)>2 and u2[-1]=='*' and u2[-2]!='.':
        # throw out trailing *
        u2 = u2[:-1]
    if '/' in u2:
        # if there is a '/' in url: relatively clear; trim, and do final checks
        u2 = u2[:u2.index('/')]
        if len(u2)>2 and u2[-1]=='*' and u2[-2]!='.':
            u2 = u2[:-1] # remove extraneous * again
    elif '*' not in u2:
        # Single domain; mostly homepages, e.g. *kreditkarte.ing-diba.de*
        # typically clear, but can include http*://*amazon* and garbage
        pass
    else:

```

```

#messiest case and where things can go wrong.
u2 = u2[:u2.index('*')] if u2[-2:]!='.*' else u2
pass

if '.' in u2:
    # require dot to avoid localnames
    # in cases of remaining wildcards, we control with publicsuffix itself in which part the wildcard
falls
    u3 = psl.get_public_suffix(u2)
    if u3 and '*' not in u3 and '.' in u3:
        return u3
    return None

n_impos,n_missing,n,cnt = 0,0,0,0
fw = open('outcome.txt','w')
fw.write('url\tdomain\n')

for s in open('ds2_urls.txt'):
    if s=='\n':
        continue
    url = s.strip()
    d_h = extract_domain(url)
    if d_h:
        simple_urls[url] = d_h
        fw.write('%s\t%s\n'%(url,d_h) )
    else:
        other_urls[url] = None
print 'total-urls: %d + %d\nmatching unidentified...' % (len(simple_urls),len(other_urls))

for url in other_urls:
    if url in (...):
        # ignore these, they mess up RE
        continue
    # convert url with wildcards into a regular expression
    regstring = ''
    for c in url:
        regstring += '.*' if c=='*' else '.' if c=='#' else '['+c+']'

    # search for regular expression matches
    matches = {} # {domain_matched: count}
    for surl,d in simple_urls.items():
        re_result = re.findall(regstring,surl)
        if re_result:
            assert len(re_result)==1
            matches[d] = matches.get(d,0)+ 1

    # Now three outcomes: no match, one or multiple matches same as before, or different from before
    d_pick = None
    if not matches:
        # no match, let's try text matching as last resort
        textmatch = 'ebay.com' if 'ebay' in url [...] else None
        d_pick = textmatch if textmatch else 'MISSING'
        n_missing += 1 if d_pick=='MISSING' else 0
    else:
        matches = sorted(matches.items(), key=lambda x: x[1], reverse=True)
        d_pick = matches[0][0]
        # one or more RE matches, and not equal to old-domain-results.
        if len(matches)>1:
            # if multiple RE matches. We chose best match only if we are more than 70% confident;
otherwise not
        p = matches[0][1]*100 / sum(v for k,v in matches)
        if p<70:
            # do text matching as last resort again
            textmatch = 'ebay.com' if 'ebay' in url [...] else None
            d_pick = textmatch if textmatch else 'IMPOSSIBLE'
            n_impos += 1 if d_pick=='IMPOSSIBLE' else 0

    n += 1
    if n%100==0:

```

```

print '.',
    fw.write('%s\t%s\n' %(url,d_pick) )
print 'wildcards-impossible: %d' % (n_impos)
print 'wildcards-missing: %d' % n_missing
print 'wildcards-updated: %d' % (len(other_urls)-n_missing-n_impos)
fw.close ()

```

The result of the above script is finally stored in a MySQL table called 'ds2_url_domain' as displayed in Figure 20. Table 'ds2_filenames_domains_urls' contains file names of the extracted domains and connects 'ds2_url_domain' table to the 'files' table.

In short, the targeted domains extracted from Zues Malware configuration files can be categorized into the following groups:

- Personal online banking
- Corporate online banking (mainly for North American small businesses)
- Investment and online trading sites
- Credit card services
- Extremely popular global websites (e.g. Amazon, eBay, Facebook, etc.)

4.3.2 INFECTION TIME

As it is discussed before, the meta-data of the actual configuration text files such as infection time is available in the 'Zeus' database in MySQL tables. In case we would be able to extract the name of the configuration text files, we would be able to find their associated infection time from 'files' table in MySQL database. The 'firstseen' column in the 'files' table in the figure below displays the time when the configuration file is first seen.

4.3.3 RC4 KEY (BOTNET KEY)

Every bot in a botnet uses a unique password (RC4 key or key-state) associated with the botnet which will be used for decrypting the configuration file(s). These rc4 keys are mostly unique for each botnet and in this research it is RC4keys are assumed as unique identifiers of the botnets. As it is clear from table 'keys' in Figure 20, each filename in 'md5' column is connected to a key (which is the rc4 key) and a keyID via the 'configs' table. In this research the keyIDs stands for botnet ids.

As we discussed in chapter, one of the primary variables that can be extracted from our dataset is the botnet key or the botnet ID.

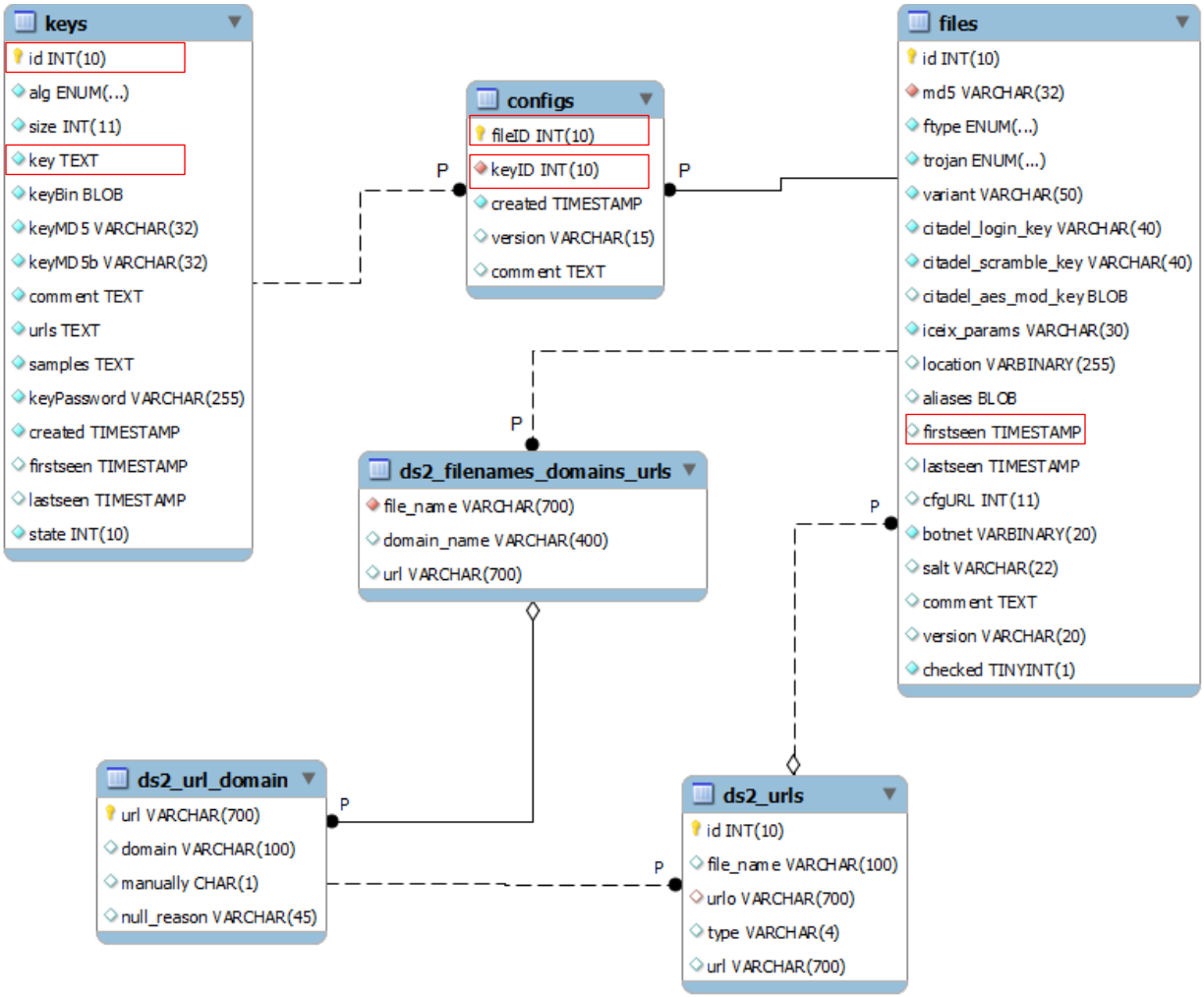


FIGURE 20- OVERVIEW OF THE MAIN TABLES USED IN THE RESEARCH

4.4 SUMMARY

In the initial sections of this chapter, the Zeus malware data that is used in this research for our empirical work is introduced and it is explained how the Zeus financial malware actually works while executing online banking fraud. The specifications of the data are also explained in details in this chapter. In order to aggregate and store the initial key information from the data we have, we extract the three main values from the Zeus malware dataset:

Value	Details
Domain name	The targeted domains
Infection Time	Time of attack
RC4 Key (Botnet key)	Botnet's unique identifier

With the information being extracted from Zeus malware data in this chapter, we now have the number, and time of the configuration files sent by each botnet in 2009 until first quarter of 2013. We used botnet's RC4 key to distinct different botnets from each other. In the next chapters of the report, this information will be used to explain the Intelligence that can be extracted from Zeus malware data.

Chapter 5 - Extracting Intelligence from Zeus Configuration Files

Introduction

In the previous chapters, we discussed that, our malware configuration data should be aggregated into a certain usable level. Thus, in the primary stage, three values, Targeted domains, Infection Time and Botnet ID are extracted from the big dataset and stored in MySQL tables.

Referring back to the objective of this research, we would like to see if any intelligence in regards to the target selection of cybercriminals could be found from Zeus malware configuration files. To investigate this, next to aggregating the data for analysis, we need to list all of the possible intelligence that can be extracted from the data using the preliminary information being extracted from the data in the previous chapter. This will ultimately help us to check whether it is possible to define any dependent variable from the information available in the Zeus malware configuration files. Accordingly, in this chapter, we will answer the following sub research question by investigating and exploring the potentials of our dataset in order to define an intermediate dependent variable:

Sub Question 3: *Is it possible to build a dependent variable to explain patterns of target selection in online banking attacks discerned from the instructions available in the Zeus malware data?*

5.1 DISCUSSIONS ON INTERPRETATIONS OF MALWARE CONFIGURATION FILES

5.1.1 CURRENT INTERPRETATIONS

Before exploring the dataset, it is crucial to mention what distinguishes this research from already existing relevant reports published by security firms by now. The reason why this research is unique and different from other existing research in online financial fraud field can mainly be attributed to the way the information in malware configuration files has been interpreted in this research comparing to the other already existing papers and security reports.

As it is presented in chapter four, from the malware configuration data it can be determined that how many times a domain is attacked. Basically, this is the number of times a domain name is seen among all of the configuration files in the period of 2009 till first quarter of 2013. These counts are called the 'raw counts' of the attacked domains which may not be necessarily the valid representation of the number of times actually a domain is attacked. Taking the example of Symantec security report published by (F-Secure, 2012a), they published a list of top 20 most attacked domains retrieved by SpyEye malware in 2012. The way they made the list is identical to the definition of the raw counts.

5.1.2 WHY ARE THE CURRENT INTERPRETATIONS INCOMPLETE?

Recalling the example of the SpyEye malware from previous section, in SpyEye malware, the configuration is built into the binary. So configurations are released as often as the binary is changed. These changes maybe in line with the changes in antivirus software because SpyEye needs to evade detection of antivirus software. These updates are extremely related to the botnet operators. This can also be true in case of other malware where a configuration file may be updated due to mistakes in the code, human faults etc.

This implies that for instance a bot-master may update a configuration file once per two days while it could take more time for a smaller botnet, a lazier bot-master, or even a botnet with stable attack code. Also, a botnet may update the configuration file less often because it uses wrapper around the binary for not being detected by antivirus, or simply because the configuration file cannot be detected by the antivirus due to problems in defense

method of antiviruses. Therefore, it can be concluded that the number of configuration files per day sent by a botnet may have little relationship with the actual attacks.

Likewise, it is probable that the result of raw counts of attacked domains in the configuration files may not represent the number of actual attacks and be over/under counted. For instance looking at the example below, the configuration files of three different botnets are displayed in a same week. With 'raw counts' we would say botnet 1 attacked 'ebay.com' three times this week and botnet 3 two times, so in total 'ebay.com' is attacked 5 times in this week. However, in practice, all of the configuration files sent by botnets in one week are only the updated versions of the previous ones and this does not indicate that the attack is performed three times in that week by botnet 1 for instance.

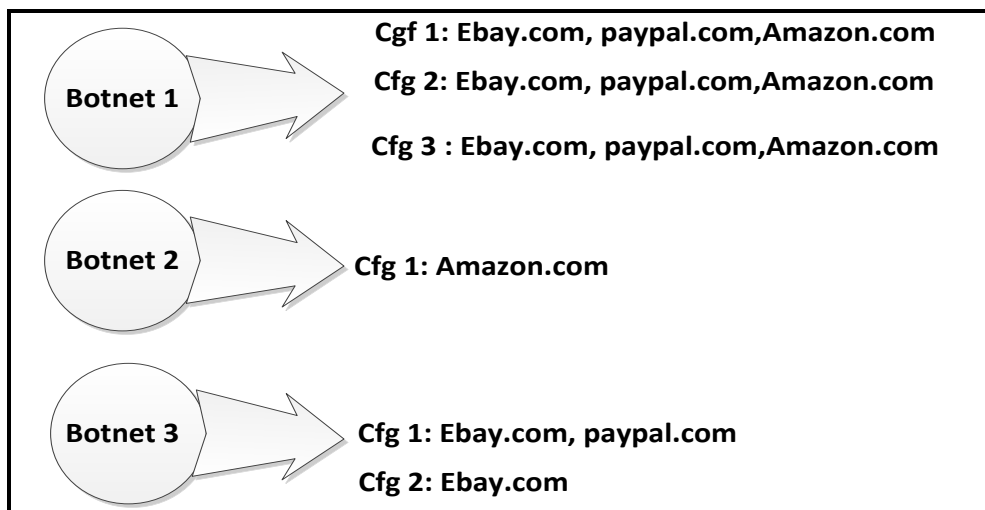


FIGURE 21- EXAMPLE FOR DISPLAYING WHY 'RAW COUNTS' ARE INCOMPLETE

5.2 LIST OF EXTRACTABLE INTELLIGENCE FROM ZEUS DATA

In the previous section, we argued that some people use raw counts to express the number of times a target is attacked. Later we argued about the reason why raw counts cannot be completely correct indicators for number of attacks. Therefore, in this section we will try to find an dependent variable that would be able to relatively express target selection by cybercriminals.

As the starting point, we will discuss all of the possible intelligence that can be extracted from the Zeus malware configuration files. Next, we will categorize our listed intelligence from the data to determine the variable that is able to explain the target selection by cybercriminal from the data in line with our main research question. In the following section we will explore a series of variables and will discuss the intelligence that they can express from the Zeus malware data as well as the limitation they may impose.

Figure 22, is built upon the conceptual framework in chapter two, and displays the relationship between the RAT three categories, intermediate variables and independent variables. In this model we use the so-called two-level theory in which basically outcomes are explained with causal variables at two levels of analysis that are systematically related to one another (Goertz & Mahoney, 2005).

The variables in the basic level contain the main causal variables and outcome variable of the theory as a whole. These variables form the building blocks of two-level theories. In our model we called them 'intermediate variables' because in they act as the connection between our main dependent variables which is 'online banking fraud' and independent variables in our conceptual model. However, at the end of this chapter we will select the

dependent variable that we intend to measure in our analysis section. The variables that are marked as independent variables in our model are actually the secondary level variables that are less central to the core dependent variables, the Zeus attack. These variables act as explanatory factors that may influence the occurrence of the core variables. In the following sections, we will explore the Zeus data to see what kind of these variables can be extracted from the Zeus data itself.

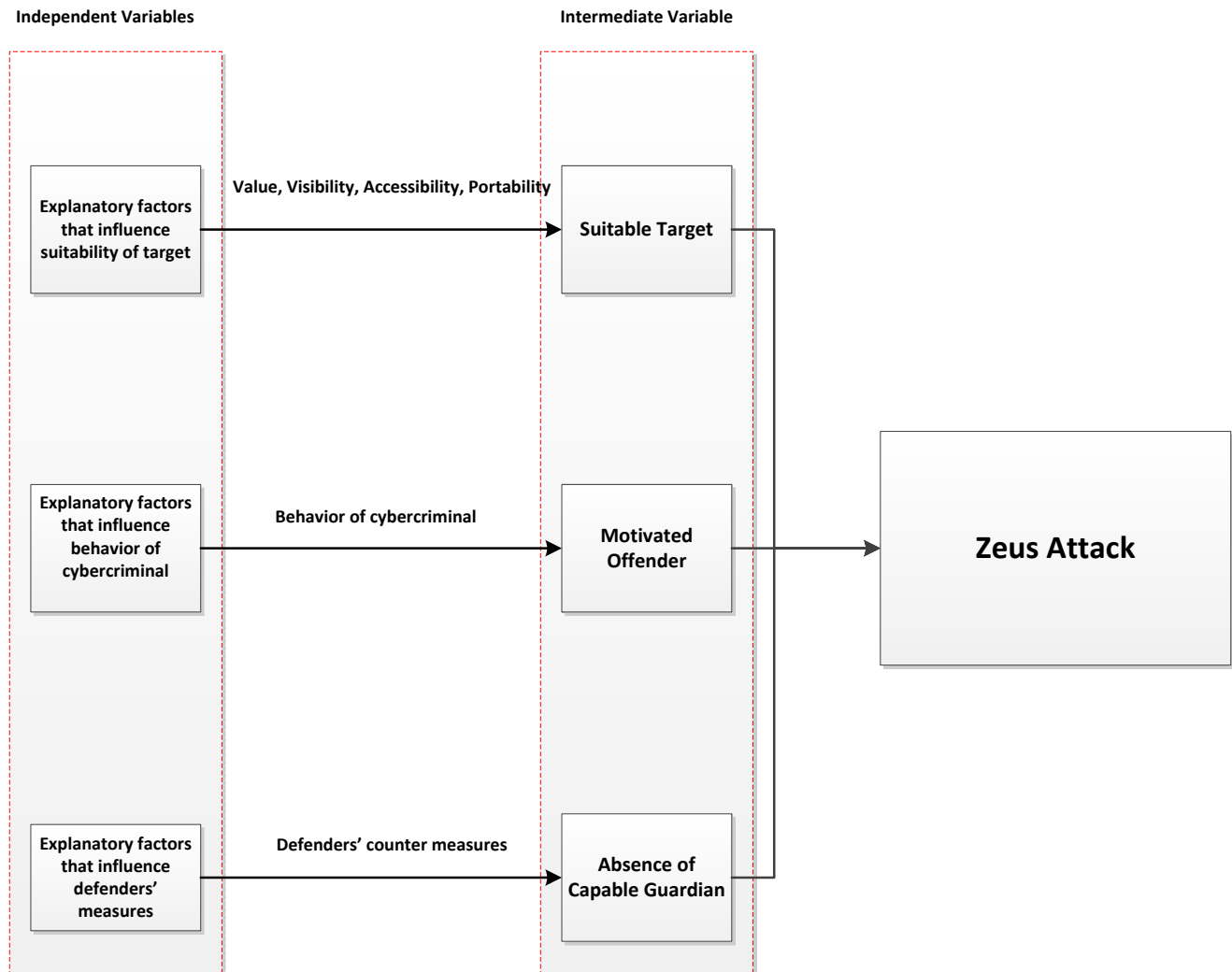


FIGURE 22- INTERMEDIATE FRAMEWORK BASED ON THE RELATIONS IN THE CONCEPTUAL FRAMEWORK

1. TOTAL NUMBER OF TIMES A DOMAIN IS REPEATED IN ALL CONFIGURATION FILES (OVER A CERTAIN PERIOD)

This variable ‘total number of times a domain is repeated in all configuration files (over a certain period)’ is identical to the raw counts that are already discussed in the proceeding section. Although we have already argued that this variable cannot be a proper indicator for target selection by cybercriminals, we decided to include this variable in our list of potential intermediate variables, to understand how different this variable is comparing to other variables. We will try to make the comparison clearer by presenting a number of visual descriptive statistics in the following section.

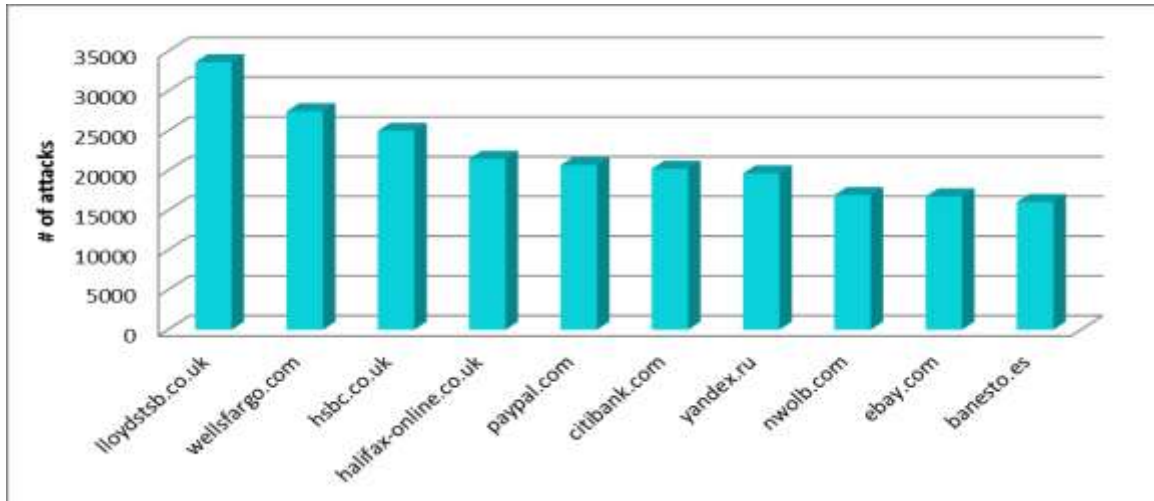


FIGURE 23- TOP 10 DOMAINS ATTACKED BY ZEUS MALWARE JAN 2009-MARCH 2013

Figure 23 displays the top-10 attacked domains according to raw counts. The rankings indicated in this figure are not be the actual indicator of the number of times these domains are attacked according to the reasons we have explained the previous chapter. Later on in this section, we will see that with another variable, these ranks are totally different in terms of number of times domains top-10 domains are attacked and their order.

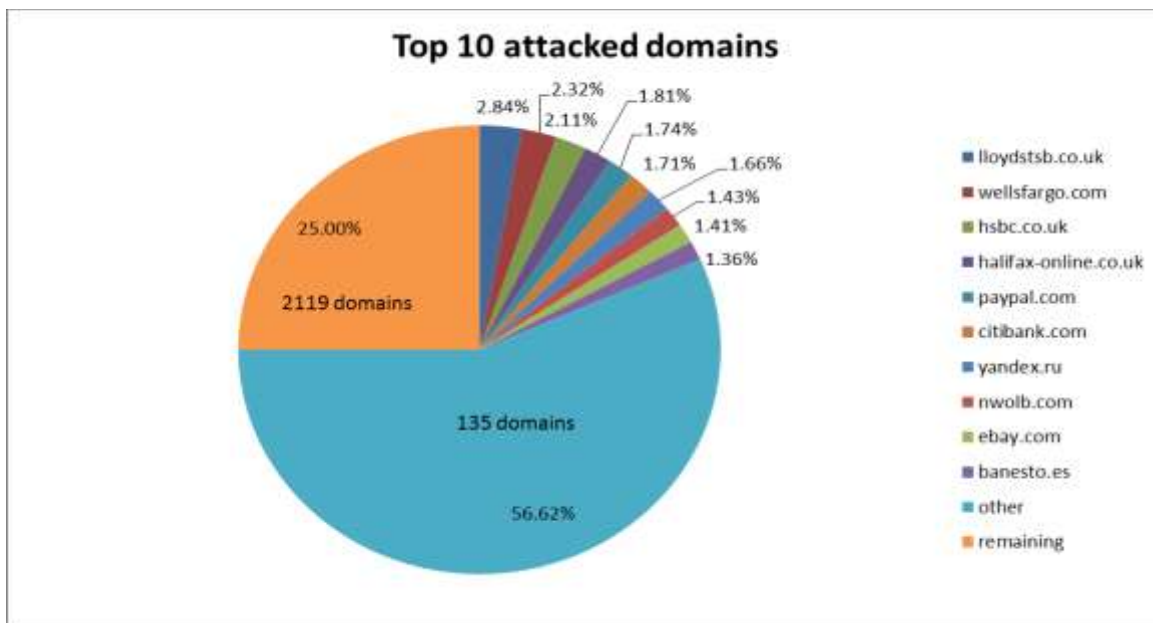


FIGURE 24- TOP 10 MOST TARGETED DOMAINS ATTACKED BY ZEUS MALWARE JAN 2009-MARCH 2013

The above pie chart (Figure 24) illustrates the percentage of top-10 attacked domains out of the overall percentage of attacked domains. As it can be seen from Figure 24, the top-10 attacked domains have approximately 20% of all attacks out of 2254 unique domains. This indicates that the 2244 domains all together have only 80% of all of the Zeus malware attacks.

In addition, as Figure 24 displays, 135 of domains account of 75% of attacks while 2119 of domains account only for 25% of attacks. This can be understood more precisely using the Figure 25. Looking to the Figure 25, it can be recognized that 10% of all domains count for 90% of all attacks. This leads us to the concept of 'power law' or 80-

20 rule. That is, a minor part of a group has dominant influence and the majority or the long tail has a little influence.

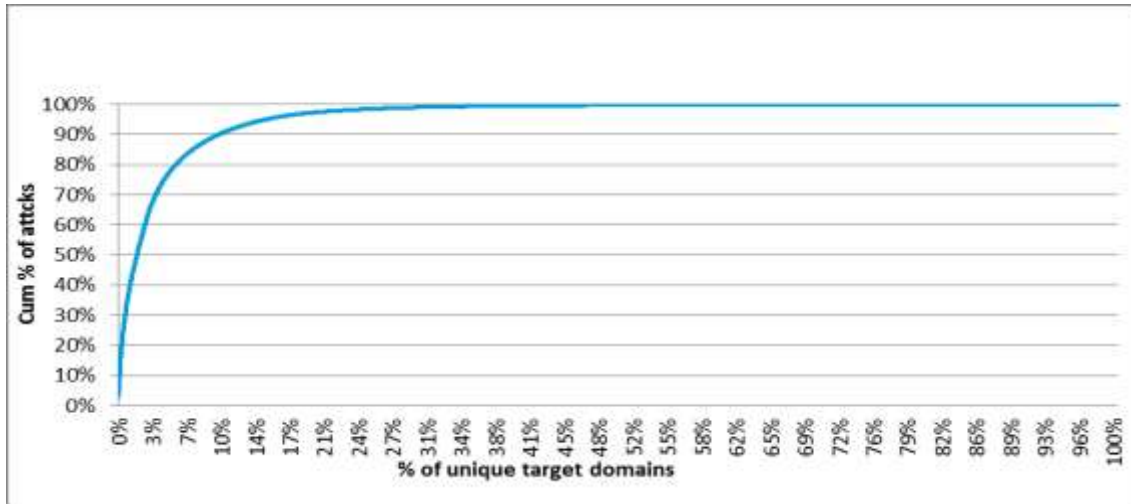


FIGURE 25- POWER LAW (CUMULATIVE PERCENTAGE OF DOMAINS COUNT FOR PERCENTAGE OF ATTACKS)

Intelligence: This variable is able to provide an insight about top attacked domains over the whole dataset.

Limitation: The reported 'number of attacks on each domain' by this variable may not 100% reliable because number of attacks may be over counted, under counted or inflated.

2. NUMBER OF WEEKS A DOMAIN IS UNDER ATTACK

'Number of weeks a domain is under attack' is a good indicator for illustrating persistence of attacks by each botnet as a different aspect of target attractiveness. We should note that this metric could be calculated in two ways. Either, the included weeks are only those when configuration files have been received by a botnet, or the weeks between the times that a configuration file is received by a botnet are also counted with the assumption that the botnet has been active in these intermediate weeks as well. The second assumption is closer to reality because as C&C servers might not often update configuration files when the attack is stable. However, in the periods that no configuration file is sent we can assume that the configuration file did not require to be updated and cybercriminals are still working with the previous version. This simply implies that previous targets in the target list of configuration file are still under attack.

TABLE 9- NUMBER OF WEEKS TOP-10 DOMAINS ARE UNDER ATTACK

Rank	Domain	# Weeks under attack
1	lloydstsb.co.uk	216
2	wellsfargo.com	216
3	hsbc.co.uk	216
4	halifax-online.co.uk	216
5	paypal.com	216
6	Citibank.com	216
7	yandex.ru	216
8	nwolb.com	216
9	Ebay.com	216
10	Banesto.es	216

Table 9 displays the number of weeks that the ten top attacked domains indicated by raw counts of the previous variable, were under attack. As it can be seen from the table, the entire top 10-attacked domain has been under attack for the same number of weeks while they have different ranks by raw counts. Because our Zeus malware data set is available for total number of 216 weeks, we can understand that these domains were under attack in the whole period of Jan 2009 until March 2013.

Figure 26 displays the number of weeks that domains were under attack. As it can be seen in this figure, a certain number of domains (about 88 domains) have been persistently under attack over almost the whole period that we have data. This is true while the average number of weeks that a domain was under attack equals to six (If we calculate the median of the number of weeks that each domain was under attack). Therefore, a number of domains are almost always the target of attack no matter of which group of cybercriminals performing the attack.

Also looking at the most left column of the figure below, we can see a number of domains (1170) that attacked less than 7 times. From the list of attacked domains, we recognize that the list include antivirus companies, phone companies and banks in a wide variety of locations (countries). From this, it can be guessed that cybercriminals perform 'trial and error' on attacking some domains. The number of attacks on these domains indicates that either cybercriminals were not successful in attacking those domains and that is why they quit selecting those domains as targets. Or because these listed domains are improved in terms of security measure over time and that is why they are not selected as targets anymore. Otherwise, simply we can concluded that some of these domains are selectively attacked few time only for a specific purpose no matter of what security measures are taken from the defender side¹⁵.

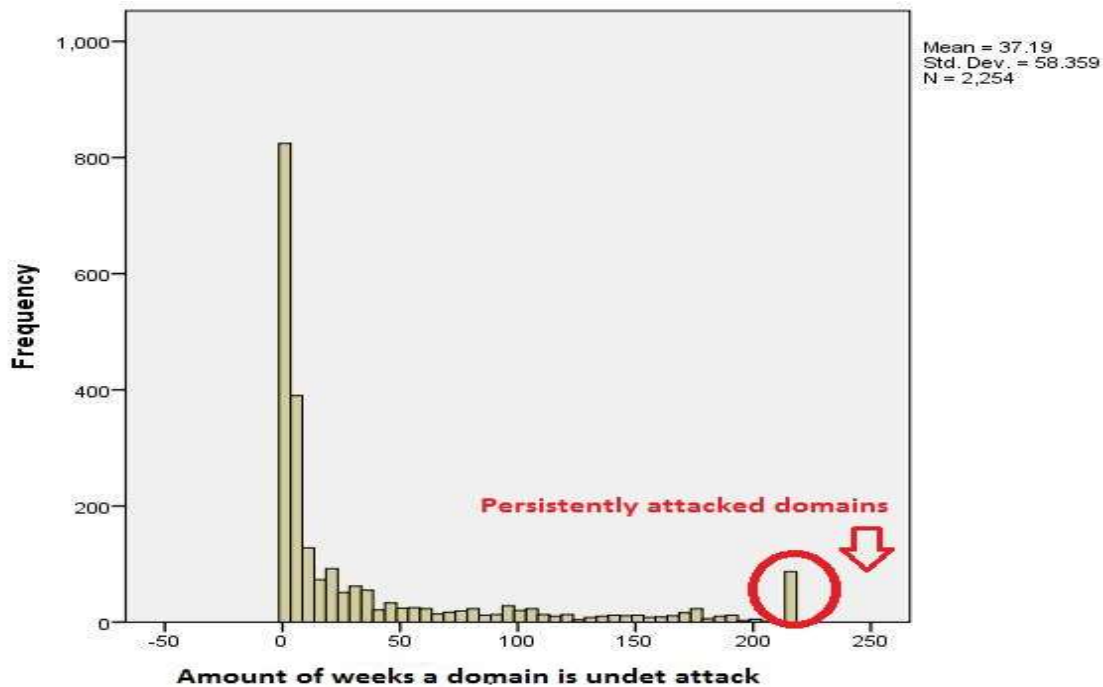


FIGURE 26-NUMBER OF WEEKS THAT DOMAINS WERE UNDER ATTACK

¹⁵ Domains indicated in the figure can also be categorized into three categories of "always attacked", "usually attacked" and "low intensity attacked" according to the number of weeks they were under attack. The more detailed characteristics of these groups can be the subject of further researches.

Intelligence: This variable is able to illustrate the attack decisions of cybercriminals looking from ‘popularity of time’ aspect. It can be concluded that a number of target domains are almost always attacked by all groups of the cybercriminal groups. It can also illustrate that cybercriminals are doing trial and error attempts on a wide variety of low intensity attacked domains (1170) and defenders are improving their security measures. Among these domains, we can find several antivirus/security software companies, phone companies, and banks. Domains that are attacked in this group (attack < 7) are from a wider variety of countries.

Limitation: it should be taken into account that, the reason behind why some of the domains attacked few times can be due to other external reasons.

3. AVERAGE NUMBER OF BOTNETS ATTACKING A DOMAIN (PER WEEK)

With variable ‘**Average number of botnets attacking a domain (per week)**’ one can determine how intense a domain is attacked among different botnets. However, this variable eliminates the limitations of raw counts (first variable) by normalizing the number of configuration files sent through considering only one configuration file by a botnet per week. The formula below displays how this variable is calculated per year. It should be noted that the variable could be calculated either per year or over the whole period that the records are available in our dataset.

$$\text{Number of attacks on a domain over a year} = \sum_{k=1}^{52} \text{botnets mentioning domain in their config file in week (k)}$$

Of course the absolute value of ‘**Average number of botnets attacking a domain (per week)**’ can be adjusted with the number of botnets active each week to create the variable ‘**Portion of botnets attacking a domain (per week)**’ as another good indicator of number of attacks on a target domain.

Looking at Figure 27, the top 10 attacked domains according to variable ‘**average number of botnets attacking a domain (per week)**’ can be seen. Comparing this figure with the one from variable ‘**total number of times a domain is repeated in all configuration files (over a certain period)**’ or the raw counts, it can easily be understood that top attacked domains and their associated ranks are different in the two variables.

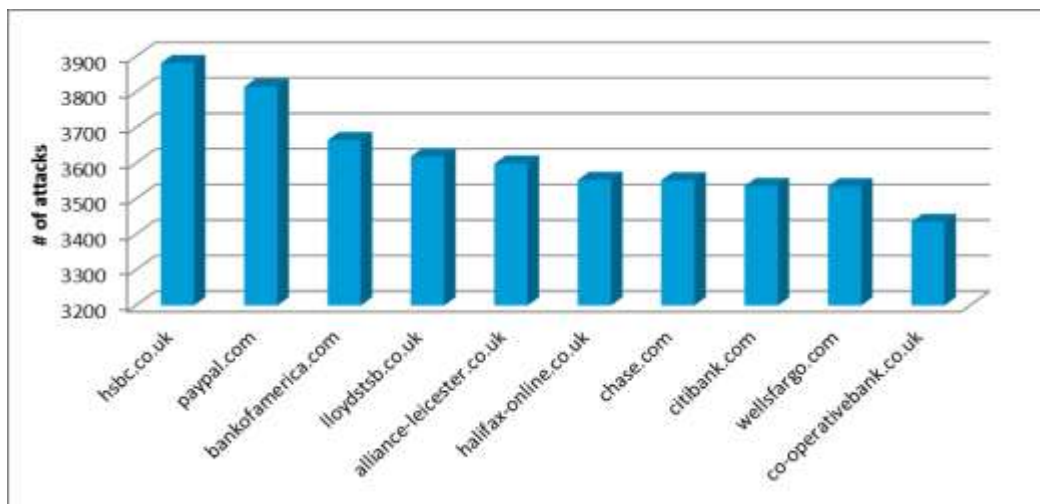


FIGURE 27- TOP 10 DOMAINS ATTACKED BY ZEUS MALWARE JAN 2009-MARCH 2013

Comparing the two pie charts, it can be seen that even the cumulative percentage of attacks on top-10 attacked domains decreases using the new variable. It can be argued that the new variable is more expressive in regards to target attractiveness because it does not have the limitations of raw counts discussed above.

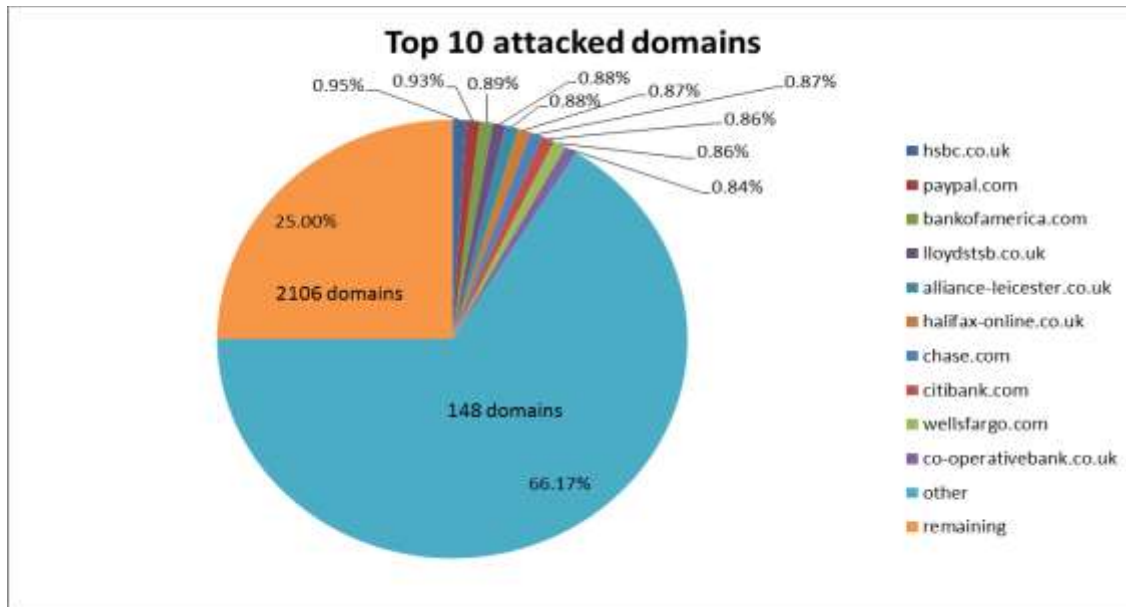


FIGURE 28- TOP 10 MOST TARGETED DOMAINS ATTACKED BY ZEUS MALWARE JAN 2009-MARCH 2013

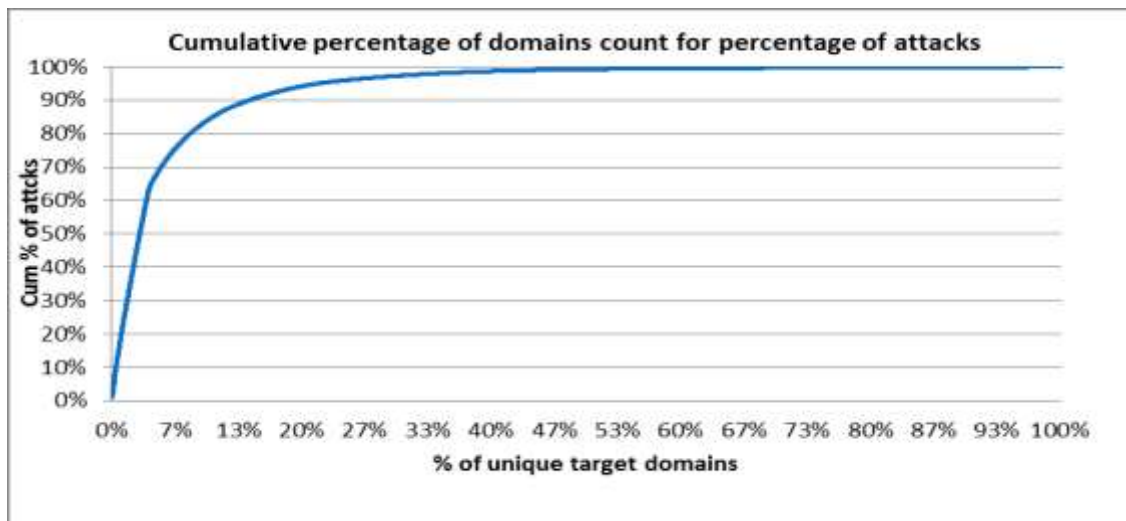


FIGURE 29- POWER LAW DISTRIBUTION (CUMULATIVE PERCENTAGE OF DOMAINS COUNT FOR PERCENTAGE OF ATTACKS)

Similar to the raw counts, still we can see the power law distribution of this variable in Figure 29. However, as it is noticeable, comparing to the raw counts (variable 1), this time 15% of attacked domains, accounts for 90% of the attacks being executed by cybercriminals which implies that there is a certain (a few) number of domains that are popular among most of the cybercriminals and a large group of domains (the long tail) that are attacked less often by cybercriminals.

Intelligence: This variable is able to illustrate number of times each botnet decided to attack different domains per week and thus lowering the repetitions in counting the number of domains.

Limitation: It may still contain a few repetitions in counting the domains that can be assumed inevitable in this stage.

4. **NUMBER OF ACTIVE BOTNETS (PER WEEK)**

Variable 'Number of active botnets (per week)' is able to illustrate the level of activity of botnets in terms of sending configuration files each week. The figure below displays the trend of botnet activity from 2009 till March 2013.0

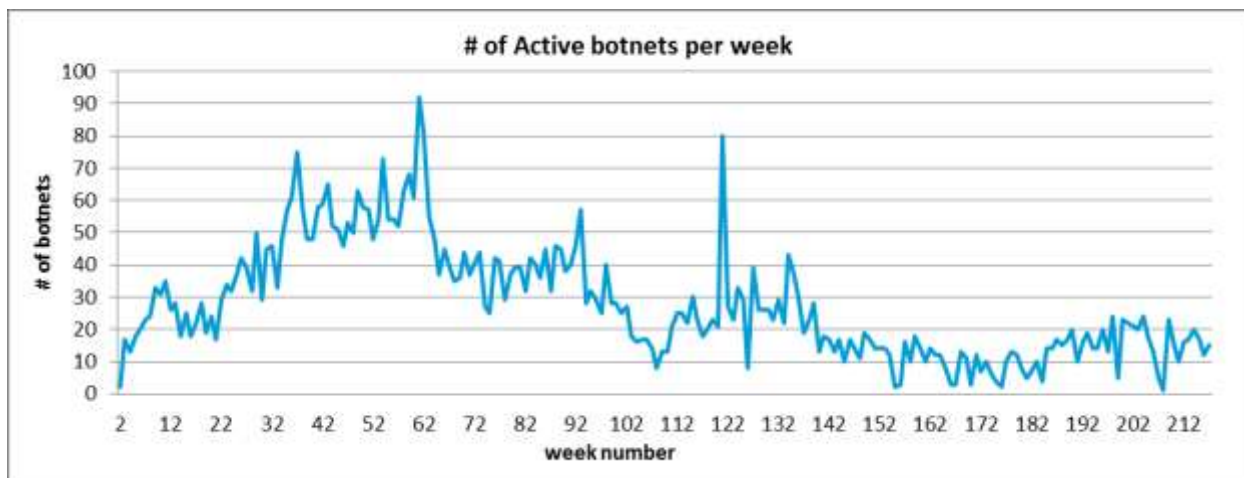


FIGURE 30- NUMBER OF ACTIVE BOTNETS PER WEEK

As it can be seen from the figure above, number of active botnets follows a decreasing trend over time. The reason behind this can be attributed to the Zeus take down efforts that are done by different governments and security firms all around the world.

Intelligence: Looking at the trend of botnet activity level each week, one can understand the decreasing trend occurs because of the Zeus take down efforts in the recent years.

5. **NUMBER OF CONFIGURATION FILES SENT PER BOTNET (PER WEEK)**

'Number of configuration files sent per botnet (per week)' is a good indicator for activity of each botnet per week over the whole period. Looking at the first graph below, the dotted line displays the number of configuration files sent each week while the solid line indicates the number of botnets have been active per week. Comparing the two lines in the following graph, it is expected that the number of configuration files sent by each botnet per week follows the same trends as the number of active botnets in that week. The reason is that the number of configuration files is calculated based on files that are sent by active botnets in different weeks.

However, it also can be seen from Figure 31 that the trend of the two lines is not identical in some of weeks. The reason can be attributed to the fact that although in some weeks the numbers of botnets have decreased, the remaining botnets were still active in terms of sending configuration files in that those weeks. That is why in some of the weeks in the graph we can see a decrease in number of active botnets when the number of configuration files sent by botnets follows the average.

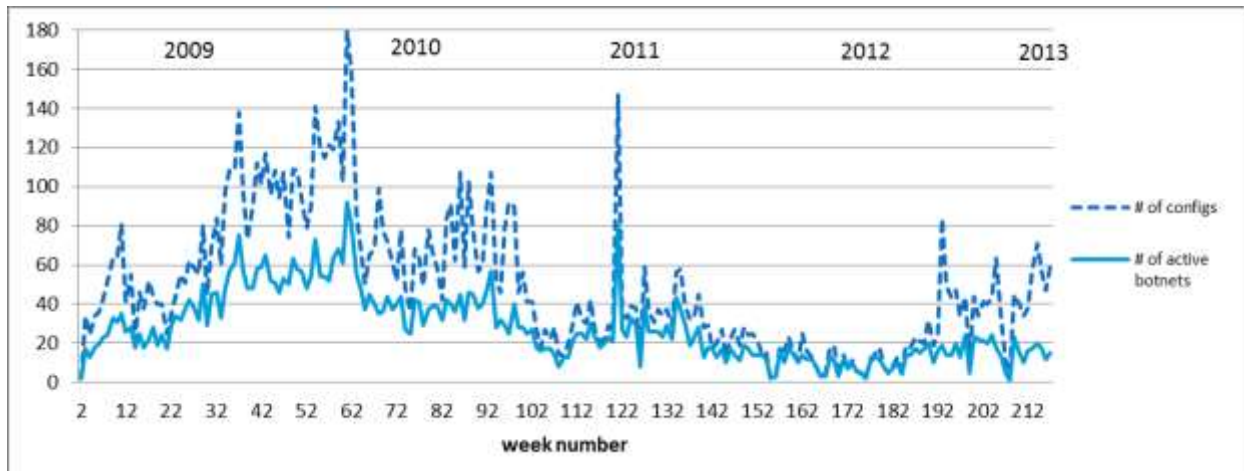


FIGURE 31- NU NUMBER OF CONFIG FILES VS. NUMBER OF ACTIVE BOTNETS PER WEEK

This comparison can be seen more clearly from Figure 32 where the number of configuration files per week is sketched taking into account number of botnets per week. The number of configuration files is divided over number of botnets per week in the following figure. As it can be seen from the second graph, the average number of configuration files sent by botnets is more or less constant during the 4 years period expect for those weeks that the number of active botnets is decreased. This implies that although fewer attackers executed the attacks in those weeks, the remaining attackers were so active in terms of sending configuration files. This can be attributed to several reasons;

- It is probable that due to improvements in the antivirus/firewall software, for executing their online attacks, attackers were forced to update their configuration files more often.
- Increase in the number of configuration files can occur due to requirement for updates as a result of human mistakes in writing the configuration files. This instance of such human mistakes have been observed several times in the inject codes available in the configuration file by the author of this report.
- It can imply that attackers are creating innovation in their attack techniques by executing trying and error.

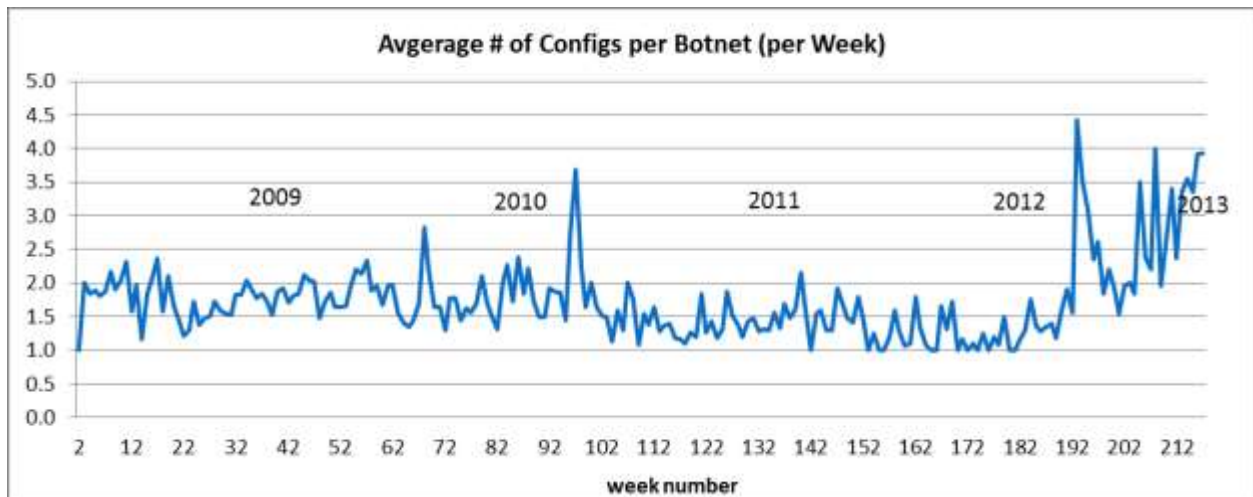


FIGURE 32- AVERAGE NUMBER OF CONFIGS PER BOTNET PER WEEK

Intelligence: By determining number of configuration files sent by botnets each week against the activity level of the botnets each week it can be seen that although after mid 2011 the number of active botnets decreased, the active botnets attacked with more number of configuration files.

Limitation: It does not provide any information about directly about target selection by these botnets.

6. NUMBER OF TARGETED DOMAINS PER BOTNET (PER WEEK)

Variable ‘number of targeted domains per botnet (per week)’ can determine the variance of the domains a botnet attacks in a specific period. Figure below displays the number of domains that are attacked each week during the 4 years period.

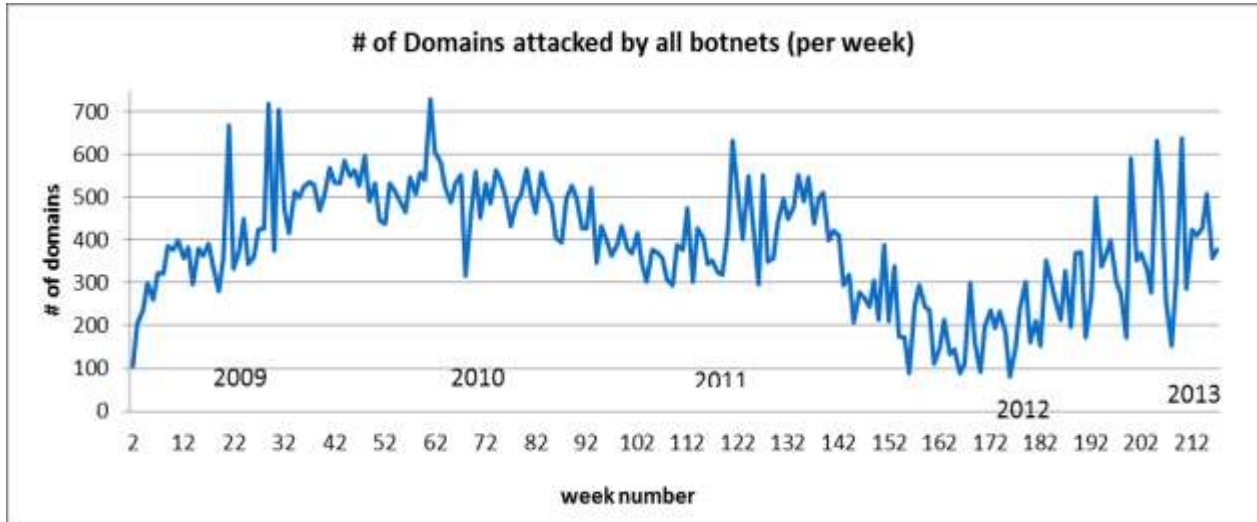


FIGURE 33- NUMBER OF DOMAINS ATTACKED BY ALL BOTNETS PER WEEK

From Figure 33 it can be recognized that the vitality of domains in the form of peaks occurs in some weeks. However, the graph can make more sense and would be more precise if we take number of botnets attacking those domains into account. The solid line in Figure 34 illustrates the average number of domains attacked per botnets per week. To be able to find the reason behind some peaks in this line, the dotted line is sketched as indicator for number of botnets active per week. From the figure, it can be understood that the peaks occur wherever the number of botnets attacked decreases. That is to say, the average of domains targeted by botnets changes rapidly in the points where the number of botnets attacking decrease. An example is January 2013 when the average has its maximum while there is only one botnet attacked 154 domains. Connecting this to the graphs presented in the previous sections about number of configuration files sent by botnets in these weeks, we can understand that the remaining botnets focused on a smaller range of targets in this period in their configuration files.

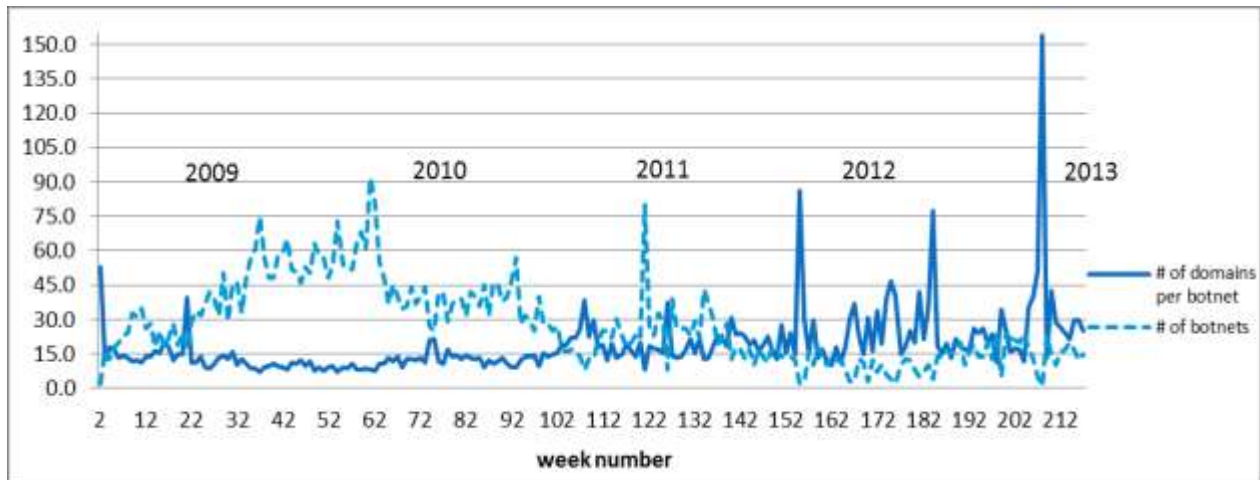


FIGURE 34 NUMBER OF DOMAINS SENT PER BOTNET PER WEEK VS. NUMBER OF BOTNETS ACTIVE

Intelligence: This variable is able to illustrate that in recent years attacks got more target-specific. This can also imply the huge target availability.

Limitation: It does not provide any information about the underlying reasons why in those weeks the variety of targeted domains are more while less number of botnets are active.

7. LIFESPAN OF BOTNETS

'Lifespan of botnets' could be an indicator for botnet activity taking into account number of configuration files a botnet sent within the period that it was active. A point of caution should be raised here that the period in which a botnet is assumed active should be explicitly specified.¹⁶ In Figure 35, a botnet assumed active during the weeks between two dates the botnet received configuration file.

¹⁶ A botnet could be assumed to be active only in the period that any configuration file is received from it. Or it may be assumed to be active during the period that no configuration file is being sent by its C&C server only because the configuration file do not need further updates. This indicates that during this period the last configuration file received by bots is being used by them and same domains are under attack. However, for using this variable, it is important that the distinction between the above points be made.

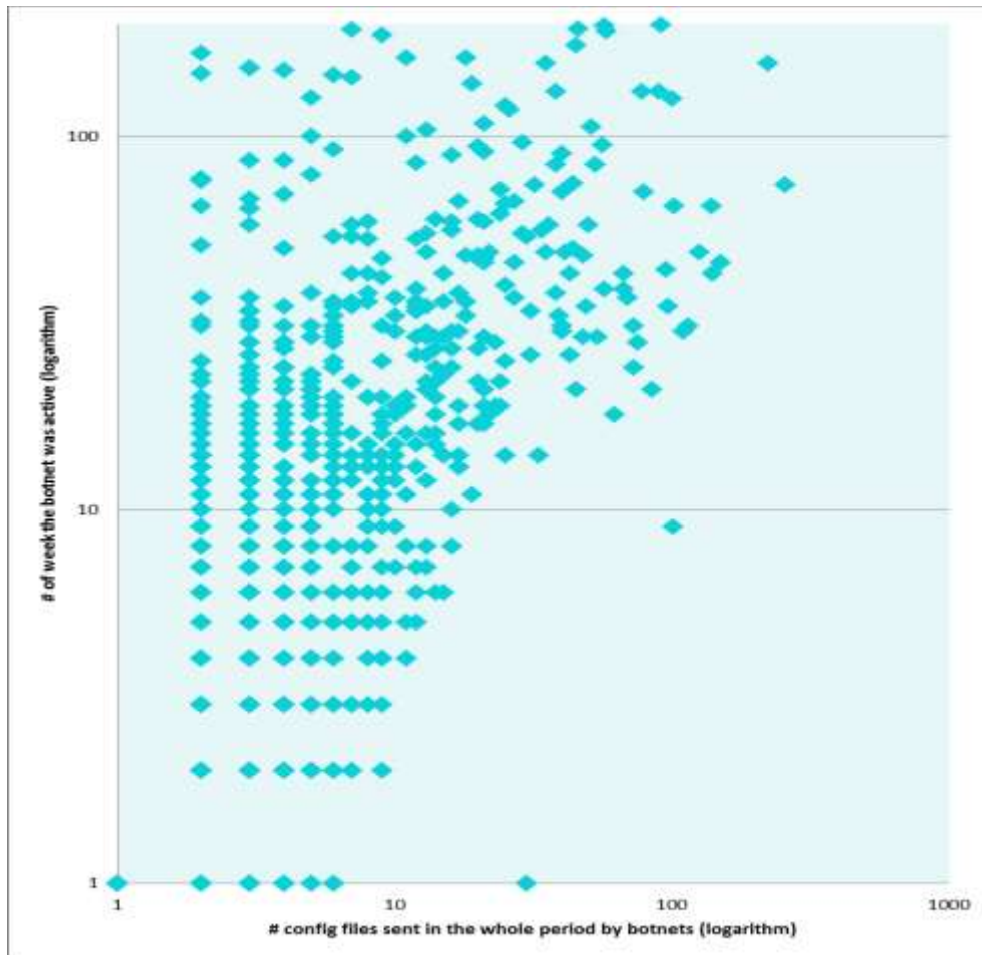


FIGURE 35- BOTNETS LIFESPAN VS. ACTIVITY (IN TERMS OF NUMBER OF CONFIG FILES THEY SENT)

Correlations

			Lifespan (weeks)	# of configs
Spearman's rho	Lifespan(weeks)	Correlation Coefficient	1.000	.906**
		Sig. (2-tailed)	.	.000
		N	2132	2132
	# of configs	Correlation Coefficient	.906**	1.000
		Sig. (2-tailed)	.000	.
		N	2132	2132

** . Correlation is significant at the 0.01 level (2-tailed).

Intelligence: There is a strong significant positive correlation between number of files that are sent by botnets and number of weeks they were active.

Limitation: It does not provide any information about specific target domains but activation period of botnets.

8. THE GEOGRAPHICAL REGIONS ATTACKED PER BOTNET

'The geographical regions covered per botnet' can determine the variance of the countries a botnet attacks in a specific period. The botnet ids are ordered by their ranks starting from the one, which attacked most. As it can be seen from the below table, diversity of countries they attack do not have a significant relation with their activity level because botnet id '802' as the most active botnet attacked 16 different countries while botnet id '2091'

attacked 27 countries having the 5th place in the table. In addition, it is interesting to mention that botnet id '841' is among the top ten active botnets by attacking only Spain and United States as targets. Moreover, it should be mentioned that a relation could be found between the two chosen countries, in terms of their languages; most of the domains registered in United States, offer 'Spanish' language as the second option in their domain web pages. Therefore, the attacker probably would not need to prepare scripts (web injects) in several different languages.

TABLE 10- INFORMATION ABOUT GEOGRAPHICAL LOCATIONS ATTACKED BY THE TOP-10 MOST ACTIVE BOTNETS

Botnet ID	# of config file sent by botnet	# of countries attacked	Attacked Locations
802	257	15	United Kingdom, France, United States ,Spain, Russia, UAE, Germany, Italy, Canada, Ireland ,Australia, India, Bulgaria, Scotland, Turkey
1415	222	10	United States, United Kingdom, Italy ,Russia , United Kingdom, Germany, Canada, Spain, Belgium,
975	150	8	Spain, Portugal ,United States, Canada ,Ireland
168	140	2	United Kingdom, United States
2091	139	25	United Kingdom, Russia, Spain, United States, Italy, Australia, France, Germany, Poland ,Czech Republic ,New Zealand ,Turkey , Antiguaand Barbuda, Ireland, China, Canada, Switzerland, Estonia, Sweden Cyprus, United Arab Emirates ,Singapore ,Portugal ,Bulgaria, India
2258	126	8	Portugal ,United States ,Italy ,Ireland, United Kingdom, Germany, Canada, Spain
841	115	2	Spain ,United States
2009	110	7	United Kingdom, United States, Spain, Germany, Ireland, Hong Kong ,Canada
1049	102	12	Turkey ,United Kingdom, United States, Australia, Italy, Spain, Canada, Germany ,Luxembourg , United Arab Emirates ,Ireland, France
1142	101	3	Spain, United Kingdom, United States

Intelligence: This variable is able to illustrate the attacking strategy of different botnets in terms of geographical locations. It can be seen that some of the most active botnets attacked a wide variety of countries while others attacked more concentrated and only to a small number of countries with similar attributes.

Limitation: It does not provide any information about specific target domains but only about popular countries and botnets.

9. PERCENTAGE OF COPYCATTED CONFIGURATION FILES PER BOTNET

'Percentage of copycatted configuration files per botnet' explains the herding behavior among cybercriminals. As it is also discussed in chapter two, behaviors may select a target or a list of targets not due to the target's specifications but only because it has been successfully attacked by other cybercriminals before. This is called 'informational cascade' which is a sub-category of herding behavior (Bikhchandani et al., 1992). Informational cascade is popular in online banking attack mainly because configuration files by different botnets are traded within cybercriminals in underground economy.

Within this research, the identical configuration files can be found using python scripts to compare the files. However, the task of comparing the configuration files and finding out the identical ones does not fit in the scope of this research due to the amount of extra of time and storage capacity it requires to be run and thus would not be performed by the author of this report. However, it can be a subject for further investigations and future research.

Intelligence: Using this variable, we will be able to determine what percent of the targets are selected regardless of their specifications and only due to copycat behavior among cybercriminals.

Limitation: Sometimes, the conclusion based on comparing targets in configuration files may not be 100% percent correct because a set of targets may only be repeated because a same attacker is performing the attack.

10. TYPE OF AUTHENTICATION METHOD DEPLOYED PER BANK (E.G. 1STEP, 2STEP)

Passing each type of authentication method requires specific steps to be taken. These steps are communicated to bots via the malware configuration file. Thus, it is possible to extract the type of authentication that the configuration file is trying to pass by analyzing the files one by one. Likewise, by comparing the web inject section of the configuration files using Python scripts; one can determine which files are passing a same type of authentication method.

Intelligence: This variable is able to provide information about the different types of authentication methods deployed by banks and attackers' reflections to these changes. We could understand why some banks are considered as suitable targets from the attackers' point of view in terms of the accessibility of target.

Limitation: It is possible that different attack methods exist for bypassing a same authentication method. Thus, determining the authentication method of each bank is not straightforward with this method especially in the cases where a large volume of malware data is under investigation.

11. CHANGES IN DEFENSE ACTIVITY LEVEL PER BANK

The defense measure taking by banks/ third party security firms working for banks may change and get more sophisticated over time. We believe that reflection of these changes could be detected in the configuration files. The reason is logical; cybercriminals will try a new attack method if their previous methods would not work anymore. Considering the fact that the configuration files contain the commands in regards the attack communicated by the C&C server, by comparing them over time for instance for some specific targets, we may be able to find relevant information about changes in defense measure employed by targets.

Intelligence: This variable is able to provide information about the changes in defense methods of banks over time and attackers' reflections to these changes determined.

Limitation: It does not provide any specific information about how the target banks are selected by cybercriminals.

5.3 CONCLUDING REMARKS

In this chapter, first, we looked at the current interpretations or in other words, Intelligence that are extracted from financial malware families in regards to the target selection by cybercriminals as a dependent variable. We argued why the current interpretations are incomplete and we build an image of how should a representative interpretation look like.

In the next step, and as our main contribution in this research, we tried to explore our case study, the Zeus malware configuration files. From the files, we found a list of possible variables that each could provide intelligence or insight about target selection by cybercriminals in online banking attacks. Some of these variables provide insight directly, such as the '**average number of botnets attacking a domain (per week)**'; others are indirectly used to normalize and build main variables such as the '**number of active botnets (per week)**'.

As the final step of this chapter, following the sub-research questioned mentioned in the introduction of the chapter, we intend to find a dependent variable based on our conceptual framework that can explain the attack decisions of cybercriminals when a target is selected. In the proceeding section, we presented various variables that provide different ways for operationalizing a dependent variable concerning how targets are selected. Among those, variable the '**average number of botnets attacking a domain (per week)**' or its normalized version '**portion of botnets attacking a domain (per week)**' are going to be used as the dependent variable in this research. The selection is mainly made for the following reasons:

- It is a good indicator of the number of attacks performed by cybercriminals because instead of counting the number of times a domain is seen in the configuration files to understand how popular it is, this metric determines the popularity of a domain by the number of botnets that attacks that domain in a week over a one-year period.
- It lowers the limitation of human errors primarily by taking into account ‘week’ as the time unit. This means that if a C&C server sends more than one configuration file to its botnet, our metric would count all as one. We argued before that logically one update of configuration file per week is rational. Secondly by taking into account the ‘number of botnets’ attack a domain, limitations resulted from raw counts are eliminated.
- The variable fits within the scope and limitations of this research.
- The required information for building the variable is available.

The next step is to find the explanatory factors that would be able to explain the occurrence of the dependent variable and build an empirical model for our analysis in the upcoming chapters. Among the variables introduced in this chapter, a number of data fields express target selection but are rather explanatory factors that can act as independent variables. Table below illustrates these independent variables extracted from Zeus configuration files based on the three categories of our conceptual framework in chapter 2.

TABLE 11-INDEPENDENT VARIABLES EXTRACTED FROM ZEUS CONFIGURATION FILES BASED ON RAT’S THREE CATEGORIES

Category	Variable
1-Profiling the Suitable Target	-
2-Profiling the Motivated Offender	<ul style="list-style-type: none"> • The Geographical regions attacked per botnet
3-Profiling the Capable Guardian	<ul style="list-style-type: none"> • Changes in defense activity level per bank

In the next chapter, we will focus on identifying more explanatory factors in regards to the target selection in online banking attacks by cybercriminals that being extracted from secondary data sources such as interview with experts of the field, and open data e-infrastructures. The focus of our investigations for explanatory factors will only be on the ‘suitable target’ category of our conceptual model due to the scope of our research. The other categories can be subject for further research in future.

Chapter 6 - Building the Empirical Model

Introduction

In the previous chapter as the core of the research, we introduced a number of variables that are able to extract intelligence from Zeus data based on the conceptual framework of chapter two. Among all of the variables introduced, the one that matched our research question the best and could act as a dependent variable for determining the number of times a domain is targeted was selected.

In this step of the research, we will enter to the stage of proof of concept that will mainly be addressed in the current chapter and the following one. Although, this part is not our main contribution in this research, we will shortly perform a series of analysis on the dependent variable that is selected at the end of chapter five. The first step for performing the analysis is to find a number of independent variables that may influence our dependent variable. In the framework presented in chapter two, some examples of independent variables were given. A few examples of independent variables that are extracted from the Zeus data have also been introduced in the previous chapter. In this chapter, we will explore the possible independent variables or so called 'explanatory factors' that can be extracted from the secondary data sources systematically, to answer the following sub question:

Sub question 4: *Is it possible to identify some explanatory factors that can explain the attack patterns and target selection by cybercriminals?*

The chapter will start by exploring the possible independent variables from the secondary data sources. Among all of the identified variables, those for which empirical data are available will be selected for the actual measurement. Accordingly, an empirical model will be built based on the relations of the independent variables and the dependent variables, which then will be converted to a set of empirical hypotheses. The chapter will be finished by explanation of each hypothesis along with its underlying rationale.

6.1 INTERVIEWING EXPERTS TO IDENTIFY INDEPENDENT VARIABLES

Normally, information on possible independent variables comes from two sources: 1. the literature, 2. interviews with field experts. Considering the fact that this research field is relatively new and that the literature is not fully developed yet, the option of interviewing experts for insights on possible independent variables is more plausible. As this research was conducted in collaboration with the Research and Development unit of Fox-IT, an opportunity existed to interview Fox-IT security experts.

In total, three security experts and one criminology researcher were interviewed. The interviews were conducted in a semi-structured format. This means that although a set of open questions were prepared for the meeting beforehand; the questions were posed as open questions to stimulate discussions with the experts.

6.2 THE IDENTIFIED INDEPENDENT VARIABLES

Table 12 lists the identified independent variables. From this section of the report onwards, we only focused on the variables in the '**Suitable Target**' category as it is in line with the scope of our research question. The below table contains variables related to the four attributes of the 'Suitable Target' category that are already introduced in chapter two namely, *portability of target*, *visibility of target*, *value of target* and *accessibility of target*.

Also as stated, the variables in Table 12 are mostly identified from interview results. It should be noted that not all of the variables will be used in our empirical model because of limited access to associated datasets and due to the limited scope of this research in terms of the amount of time and effort that is needed to gather such data.

TABLE 12- LIST OF INDEPENDENT VARIABLES BASED ON INTERVIEW RESULTS, GROUPED USING RAT CATEGORIES

Influencing Factor	Variable	Description
Portability of Target	Mule recruitment ¹⁷	Defined as the process of recruiting money mules ¹⁸ for transferring money from a victim's account to an account controlled by the money mule. The process is different in different cybercriminal groups and is dependent on a country's banking policy. This determines how fast and easily the stolen money can be reached by cybercriminals.
	Bank's cross-national money transfer policy	A bank's policy and terms and conditions for cross-national money transferring. This can influence the amount and speed that the stolen money can be transferred from a bank in one country to a bank in another country.
	Bank's online transaction clearance time	The time that takes in case of an online transaction to be completed. This is different in different banks and countries and can influence how fast and easy the money can be transferred to the cybercriminals account.
Visibility of Target	Banks with the word 'bank' in their domain names	Banks with a domain name that contains the word 'bank'. Experts believe that this could increase the visibility of the target bank for cybercriminals when they are selecting their targets from online sources and search engines.
	Herding / information cascade	Security experts argue that herding behavior exists in online banking fraud. This means that a bank or a list of banks may be targeted by a group of cybercriminals only because it was attacked previously successfully by another group, regardless of characteristics of the target itself.
Value of Target	Number of online users of a bank (heuristics: bank size)	The number of online users of a bank can influence how valuable a bank is as a target, meaning a larger pool of potential targets is available. (Bank size can be a good proxy for this variable)
	Number of Internet or fixed broadband users	As the number of people using Internet in a specific country increases, the probability of financial gain for cybercriminals also increases because the number of potential victims increases.
	The amount of time people spend online	For an online banking attack to be successful, not only should people be connected to Internet, but also they should spend time online. The time spent online is another factor that can influence the value of a target (country-specific) and thus the attackers' decisions in regards to target selection.
	Rate of online banking penetration	Defined as the amount that people do banking activities via online banking channels. As the rate of online banking in a country goes higher, banks located in that country become more valuable targets for cybercriminals as fraud targets.
	Rate of online shopping penetration	The popularity of online shopping in a country can also be used as an influential factor for selecting targets by cybercriminals. This is because it is more probable that people who do online shopping via their debit/credit card be victim of online banking fraud rather than those do not do any online shopping.
	Country's financial situation	We can argue that, it is mostly probable that banks in financially well-doing countries are more valuable targets for online banking fraud comparing to other countries. (GDP can be a good proxy of wealth of a country)

¹⁷ This variable is also mentioned in Anderson, R., Barton, C., Böhme, R., Clayton, R., van Eeten, M. J. G., Levi, M., Moore, T., & Savage, S. 2012. Measuring the cost of cybercrime.

¹⁸ Money mules are individuals recruited wittingly and often unwittingly by criminals, to facilitate illegal funds transfers from bank accounts.

Accessibility of Target	Bank authentication method	Mechanisms with which the users connect to their online banking platform. Experts believe that a bank's authentication mechanism has a decisive role in determining access to a user's online account.
	Language of the banking webpage etc.	Security experts believe that the language of banks' webpage plays an important role in increasing access for cybercriminals So for instance, due to language barriers, banks with English web pages will be more accessible to cybercriminals in terms of writing injection scripts.
	Users' online awareness	Defined as the amount of technical knowledge users have with regards to the online environment and its threats. As this knowledge increases by users doing online banking, the probability of those users getting hacked should decrease. ('Human development index' can be a proxy for this variable).
	Rate of use of firewall/antivirus products	The number of online users using firewall or antivirus products in a country. When this rate is higher, users have more security measures on their personal computer and so the probability of them getting hacked becomes lower. (Bot infection rate can be a proxy for this variable)

6.3 THE EMPIRICAL MODEL

The empirical model will not contain the entire dependent and independent variables discussed, but just the ones that gathering data for them would be possible taking into account the scope of this research. Based on the empirical model, a series of empirical hypotheses will be built to test the relations indicated in the model.

Figure 36 displays the empirical model. The model is built as follows: the conceptual framework is taken as the starting point; for the dependent variable, the metric '**average number of botnets attacking a domain per week per year**', as motivated in the previous chapter, is chosen among the possibilities. Likewise among all of the independent variables introduced in the preceding section, those that we found empirical data for are chosen.

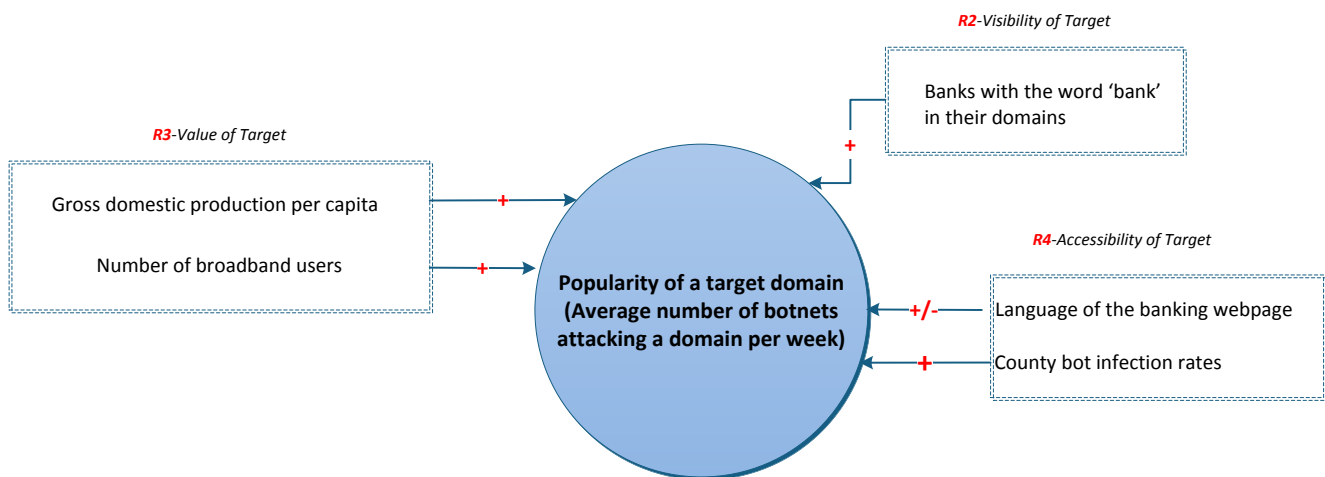


FIGURE 36- THE EMPIRICAL MODEL

6.4 PROPOSED HYPOTHESES

Based on the presented empirical model and in light of our research question, we introduce the following set of hypotheses. For each hypothesis the rationale behind it is reiterated.

Hypothesis 1: Banks with word 'bank' in their domain names are more popular among cybercriminals.

Rationale: Banks with word 'bank' in their domain name are more visible to cybercriminals than other banks because without a lot of effort it is obvious that a certain domain is a bank web page. Whereas for domains without word 'bank' in their domain names, cybercriminal requires to perform a web search, which could even get harder in the cases in which the domain name is not chosen in English but in a foreign language. So banks with bank 'word' could be noticed and recognized by cybercriminals sooner and easier and so they are more probable to be more popular as attack targets.

Hypothesis 2: A bank is more popular among cybercriminals as a target when it offers online banking in English.

Rationale: Writing and injecting English scripts requires fewer attempts (knowledge/timewise) for attackers, compared to other languages (e.g., Dutch, German or French) for which people with local language skills are required. That is why experts believe that domains that have webpages in English are more popular targets than those with pages in only local languages.

Hypothesis 3: Banks in countries with higher broadband penetration are more popular targets among cybercriminals.

Rationale: Because the probability of using online banking increases as people as more people come online and in particular when they have high-speed access. This ultimately increases the chance of gaining more value from targets from the cybercriminals' perspective.

Hypothesis 4: Banks in countries with higher GDP are more popular attack targets among cybercriminals.

Rationale: GDP can be a good proxy of wealth of a country. It is mostly probable that banks in financially well-doing countries are more valuable targets for online banking fraud comparing to other countries.

Hypothesis 5: Banks in countries with higher rate of infection are more popular attack targets among cybercriminals.

Rationale: When this rate is higher, users have less security measures on their personal computer and so the probability of them getting hacked becomes higher.

Introduction

In the previous chapters, we proposed methods for extracting intelligence from the Zeus financial malware dataset, and defined a dependent variable that fits our research question best. We also built an empirical model based on a selection of independent variables to address our sub questions, and finally, proposed a set of hypotheses. In this chapter, we will analyze our dataset using three methods: descriptive analysis, bivariate analysis, and decision-tree analysis. The details of the employed methods will be explained in the upcoming sections. Although we have previously explained the general process through which the data was prepared, for each of the methods we will briefly explain the details of the dataset created specifically for that method.

7.1 DESCRIPTIVE ANALYSIS

In chapter five of this report, we discussed a whole array of variables capable of extracting intelligence from the Zeus malware dataset. The variables in chapter five were categorized using routine activity theory into three groups. Although we ultimately chose a dependent variable for our measurements from the 'suitable target' category, we explained that other variables in this category and in the 'motivated offender' category are also helpful and can give insight on criminal behavior and target selection. In the following section, we will discuss those variables one by one in more details concerning intelligence they can provide.

7.2 POPULATION OF THE TARGETED DOMAINS

Figure 37 displays the top ten attacked domains, calculated using variable '**average number of botnets attacking a domain**'. The pie chart (Figure 38) illustrates the percentage of attacks on the top domains out of the overall number of attacks. As can be seen from this figure, the top 10 attacked domains, equal to 0.004% of the domains in the set, host approximately 10% of all attacks. If we increase this to the top 158 attacked domains (7% of all the domains in our set), we account for 75% of all attacks. Figure 39 illustrates this in another manner and shows a power law distribution where 10% of the domains account for 90% of the attacks.

Looking into Figure 39, we can understand that firstly, there is a list of approximately 150 popular domains that capture a majority of all attacks. Secondly, the population of target domains are spread out to a certain level that even the top-10 targets captures only a small percentage of attacks. Although this is a power law distribution, it is not as extreme as what we often see in other areas of cybersecurity. For instance, in case of spam, from total number of 40,000 ASN¹⁹ 100 of them accounts for almost 50% of all attacks which is also an instance of a power law but more concentrated (Van Eeten, Bauer, Asgharia, Tabatabaie, & Rand, 2010).

¹⁹ Autonomous System Number (ASN) and in this article is used as a proxy for ISP.

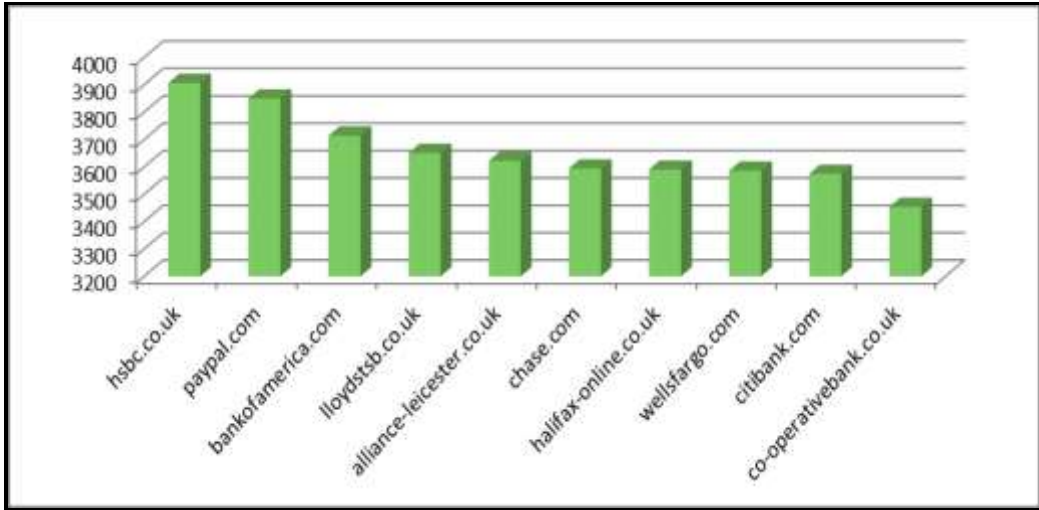


FIGURE 37-TOP-10 ATTACKED DOMAINS BY ZEUS FINANCIAL MALWARE

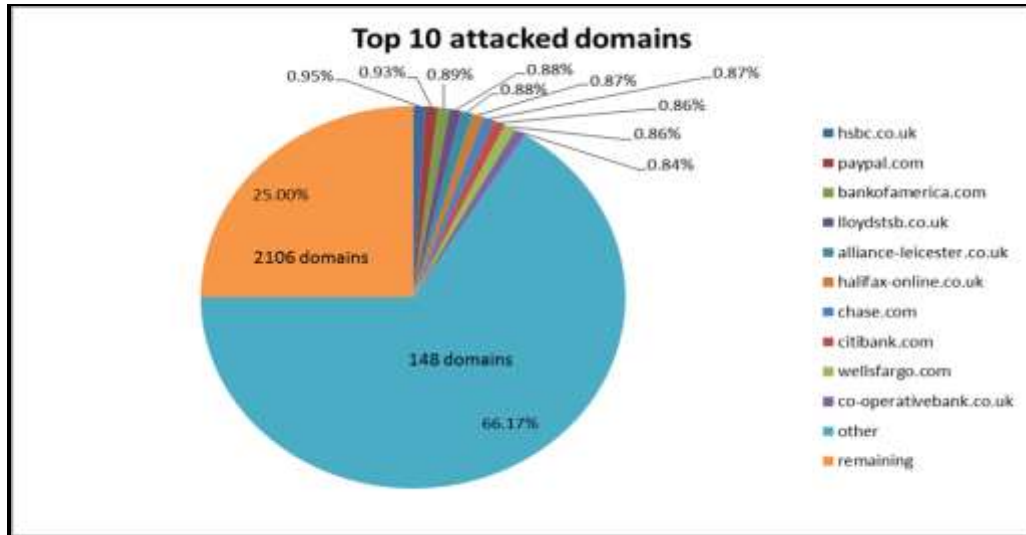


FIGURE 38- TOP-10 ATTACKED DOMAINS BY ZEUS FINANCIAL MALWARE

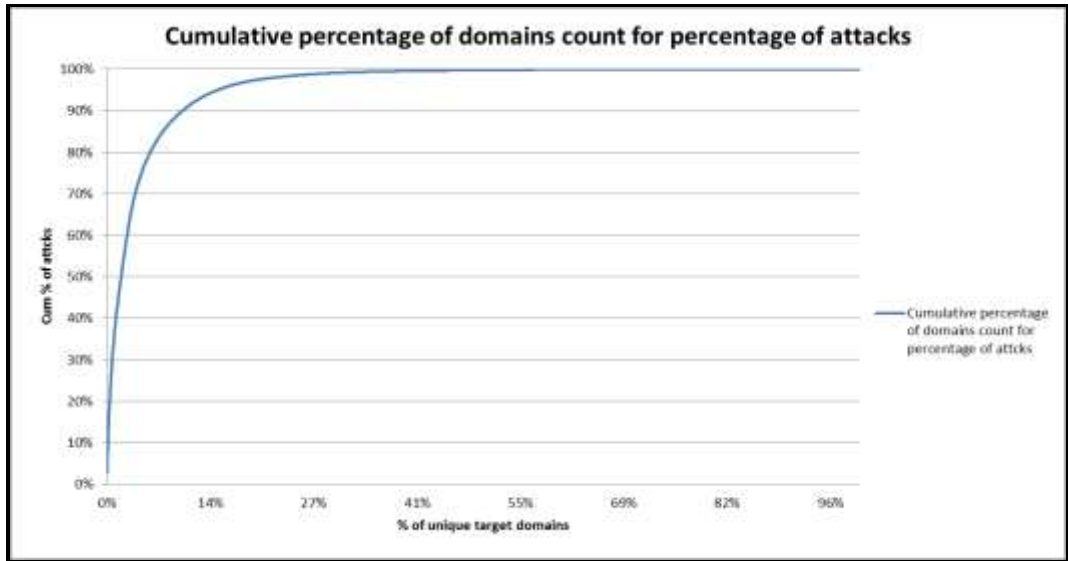


FIGURE 39- DISTRIBUTION OF THE TOP-10 ATTACKED DOMAINS BY ZEUS FINANCIAL MALWARE

7.3 BOTNET ACTIVITY TRENDS

Figure 40 displays the trend of botnet activity from January 2009 up to March 2013. The number of active botnets follows a decreasing trend over time. The reason behind this might be attributed to the Zeus take down efforts that were coordinated by different governments and security firms all around the world²⁰. It has to be taken into account that in this research each unique login key is assumed a different botnet; in practice however, a botnet may use different login keys over time. This should be marked as a limitation of our analysis.



FIGURE 40- NUMBER OF ACTIVE BOTNETS PER WEEK

The dotted line in Figure 41 displays the number of configuration files sent each week by all botnets together. The solid line indicates the number of botnets that were active in that week. Comparing the two lines, we see that they follow a similar trend. This is to be expected, as the number of active botnets is determined by whether they have

²⁰ Microsoft, Financial Groups Execute Takedown of Zeus Botnet Servers: blogs.technet.com/b/microsoft_blog/archive/2012/03/25/microsoft-and-financial-services-industry-leaders-target-cybercriminal-operations-from-zeus-botnets.aspx

sent a file or not that week. However, it also can be seen from the figure that the trend of the two lines is not identical in all weeks. In those weeks, a number of the botnets were more active in terms of sending configuration files.

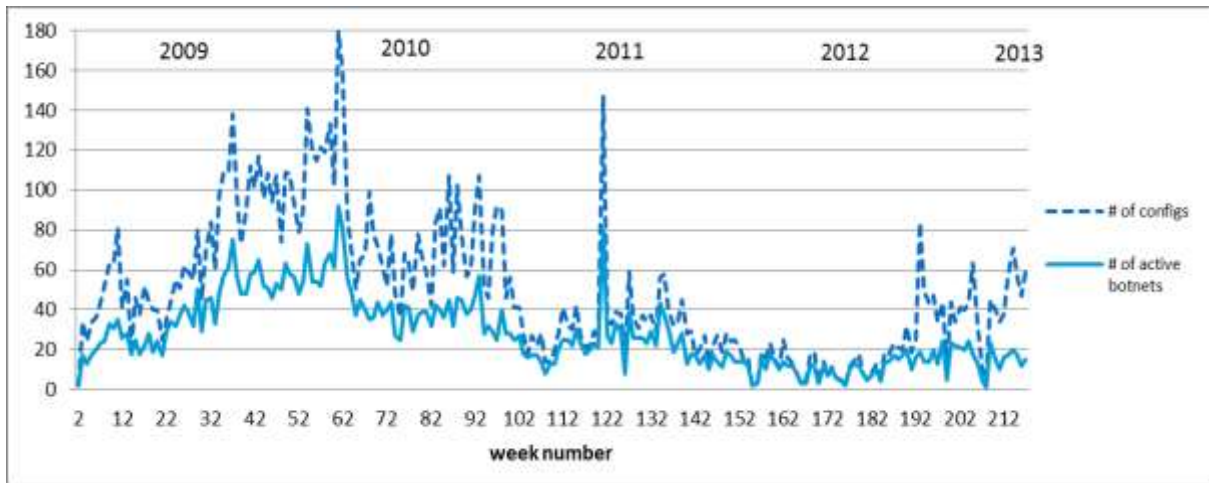


FIGURE 41- NUMBER OF CONFIG FILES VS NUMBER OF ACTIVE BOTNETS PER WEEK

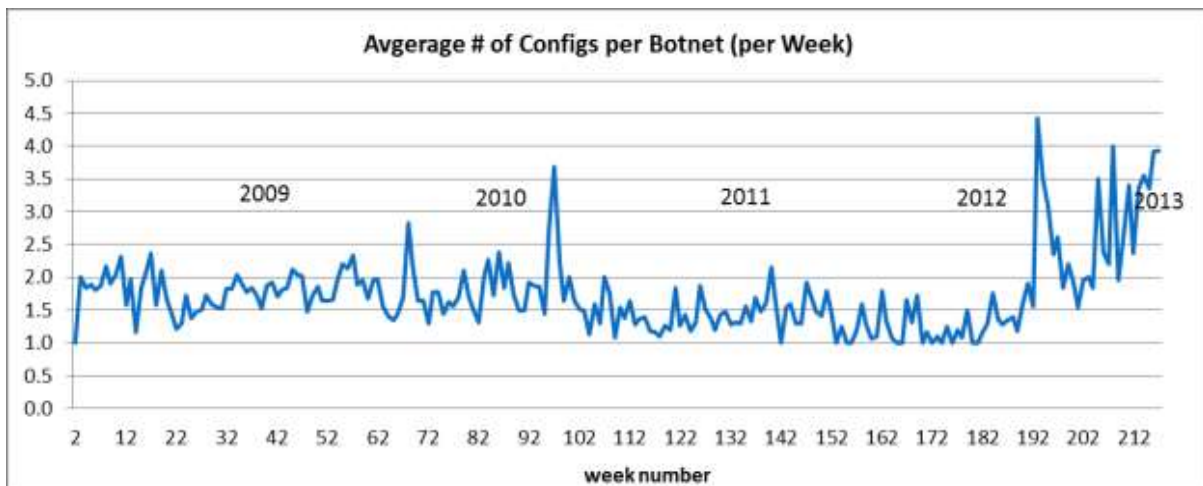


FIGURE 42- AVERAGE NUMBER OF CONFIG FILES PER BOTNET PER WEEK

If we divide the number of configuration files per week by the number of botnets active in the same week, we get to Figure 42. The peaks in late 2012 and early 2013 occur due to a decrease in number of active botnets. While still the number of configuration files sent by them is close to the average of files sent in the whole period.

Figure 43 shows a number of peaks in the average number of domains attacked by botnets in late 2012. Taking into account the decrease in number of active botnets and increase in the number of configuration files, we can conclude that the number of not-very-active botnets has gone down; the remaining are very active and they attacked a wider variety of domains, so the overall domain-per-botnet goes up. This is in line with some of the reports published by security firms such as TrendMicro (2013) that financial attacks are getting more selective and target-specific.

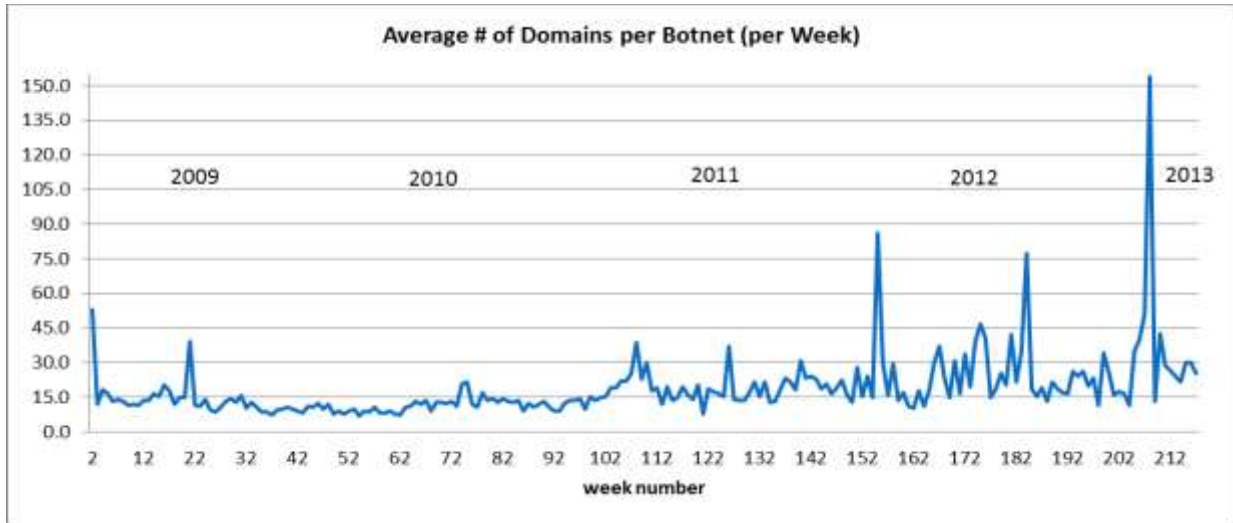


FIGURE 43- AVERAGE NUMBER OF DOMAINS ATTACKED PER BOTNET PER WEEK

7.4 ATTACK PERSISTENCE

Figure 44 displays the number of weeks different domains were under attack between January 2009 and March 2013. A small number of domains (about 90) were under attack for 216 weeks, the whole period our dataset covers. (We call these the always-attacked domains). This group of domains is usually the target of attack, no matter which criminal groups appear to be active at that point. We also have a large group of domains (about 1273) that come under attack for under ten weeks (rarely-attacked domains). Finally, there is a group of 890 domains between the two extremes that were attacked between 10 to 200 weeks (usually-attacked domains).

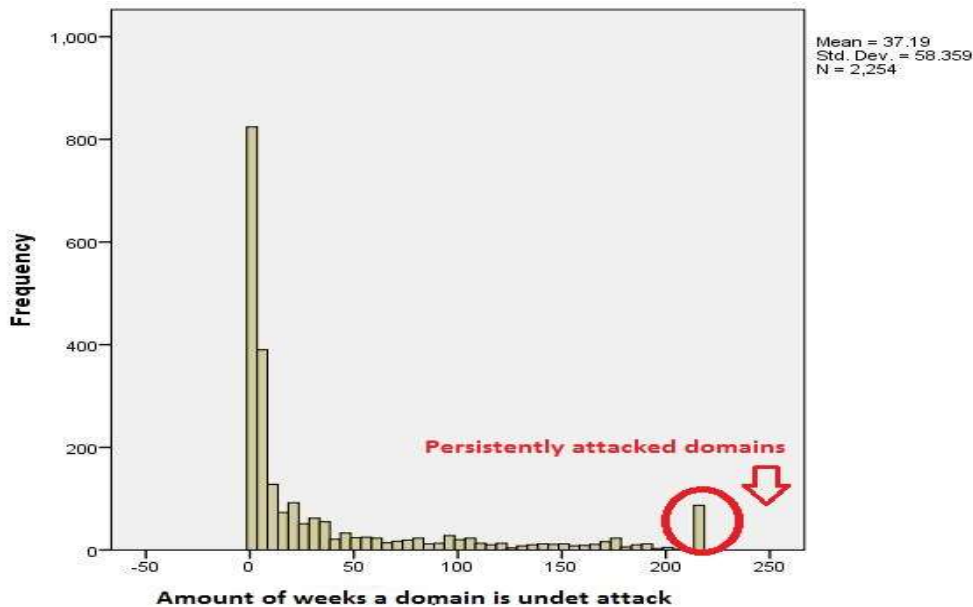


FIGURE 44- AMOUNT OF WEEKS DOMAINS WERE UNDER ATTACK

If we take a closer look at the list of **'always-attacked'** domains, we see they are mostly located in *Italy, Spain, Russia, United States* and *United Kingdom*, and a few domains from *France, Australia, Germany* and *Ireland*. That is in 9 countries out of the total 61 countries in the dataset. Additionally, 56 out of the 88 domains offer 'English'

banking webpages. The majority are banks, with an additional small group of online payment websites such as PayPal and E-Bay. The **'rarely-attacked'** domains include antivirus companies, phone companies and banks in a wider variety of countries and languages.

Figure 44 also indicates that there is a wide range of potential targets for cybercriminals. The fact that some of the targets are quite short-lived might imply that cybercriminals perform a lot of trial & error and experiment with attacking many domains; (we might even phrase this as attacker R&D). To elaborate: one could argue that either attacks executed on the rarely attacked domains are short-lived because they are not successful and the cybercriminals move on; or alternatively, that these domains are selectively attacked a few times each for a specific purpose, i.e. as part of targeted attacks.

If we compare the rank of our attacked domains (from the always-attacked to the rarely-attacked) with Alexa (2013) rankings of website popularity, we expect to see a correlation between our highly attacked domains and the popular websites by Alexa ranking.

The scatter plot and test results are displayed below. First, there is weak and significant negative correlation, which is to be expected. This implies that the domains that are longer attacked are the more popular domains in the Alexa ranking²¹ as well. Moreover, as it can be seen in the scatter plot, the **'always-attacked'** domains are mostly placed in the lower ranks of Alexa (except for some outliers). Most of the **'usually-attacked'** domains are also placed in the low to medium ranks in Alexa. The only category that contains a lot of domains in the high ranks of Alexa ranking is the **'rarely-attacked'** group of domains.

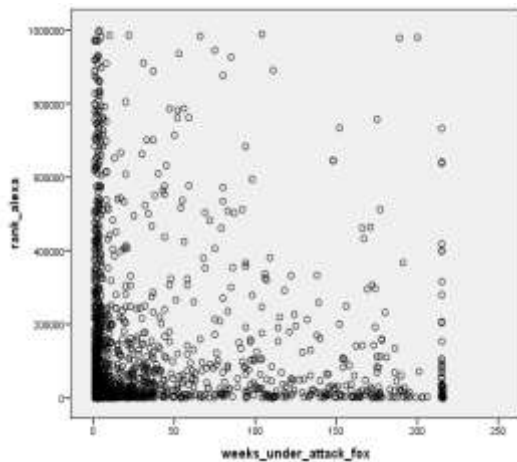


FIGURE 45- ALEXA RANKS VS. ZEUS RANKS

TABLE 13- SPEARMAN'S CORRELATION TEST BETWEEN ALEXA RANKS AND ZEUS RANKS OF ATTACKED DOMAINS

Correlations			weeks_ under_attack_ fox	rank_alexa
Spearman's rho	weeks_under_attack_fox	Correlation Coefficient	1.000	-.156**
		Sig. (2-tailed)	.	.000
		N	1492	1492
	rank_alexa	Correlation Coefficient	-.156**	1.000
		Sig. (2-tailed)	.000	.
		N	1492	1492

** . Correlation is significant at the 0.01 level (2-tailed).

²¹ In Alexa ranking the most visited page ranks 1, so the Alexa ranks are opposite of our ranks.

7.5 BIVARIATE STATISTICAL ANALYSIS

We have five hypotheses to test (relisted in Table 14); all of them will be tested using ‘**Bivariate Analysis**’ techniques. The **SPSS** software package is used for performing the analysis. For hypothesis 1 and 2 the ‘comparison of samples means’ test will be used because two distinguished groups can be identified. The other three hypotheses will be tested using the ‘measures of association’ test.

TABLE 14- LIST OF HYPOTHESES

#	Hypothesis
1	Banks with word ‘bank’ in their domain names are more popular among cybercriminals.
2	A bank is more popular among cybercriminals as a target when it offers online banking in English.
3	Banks in countries with higher broadband penetration are more popular targets among cybercriminals.
4	Banks in countries with higher GDP are more popular attack targets among cybercriminals.
5	Banks in countries with higher rate of infection are more popular attack targets among cybercriminals.

7.5.1 SPECIFICATIONS OF THE DATASET

The final dataset that is used in the statistical analysis has a number of specific properties, which are going to be addressed in this section.

7.5.2 LEVEL OF ANALYSIS

The unit of analysis is different in the following hypothesis. For hypotheses 1 and 2 the unit of analysis is the ‘domain level’; the rest of the hypotheses are tested on the ‘country level’. This relates to whether the independent variable used in the hypothesis is domain level or country level.

7.5.3 DETERMINING LOCATION & LANGUAGE INFORMATION

Our dataset contains a wide variety of targeted domains, spanning 61 countries from across the world. The location of each attacked domain was determined manually, with the help of a student assistant, by visiting each domain. This is especially important for generic TLDs. Out of the total 2255 domains in the set, we looked at all those appearing more than 10 times in the configuration files in the whole period, which equates to around half the domains. The language of each domain was also checked in this process.

7.5.4 SELECTING GEOGRAPHICAL REGION

Given the nature of bivariate analysis, a correlation could hold for a group of specific countries, but not for all the attacked countries; in order not to lose the potential patterns in our dataset, we will perform our tests on two groups: countries in the EU region (listed in Table 15) and all regions. A point of caution should be raised that we do not include countries that are not in our dataset. An alternative approach could be to include such the data about the countries that are not in our dataset in the analysis with zero attack counts for its relevant period²².

TABLE 15- LIST OF EU COUNTRIES IN/OUT OF OUR SAMPLE

Region	Country
EU in our sample	Austria, Belgium, Bulgaria, Cyprus, Czech Republic, Estonia, Finland, France, Germany, Hungary, Ireland, Italy, Latvia, Lithuania, Luxembourg, Netherlands, Norway, Poland, Portugal, Romania, Spain, Sweden, United Kingdom
EU not in our sample	Croatia, Denmark, Greece (no data on it), Malta, Slovakia, Slovenia,

²²In our empirical work, we tried this approach and included the EU countries attacked for 2012. This however did not have a significant impact on the result of the bivariate tests; so we did not expand it to the whole dataset. This approach can be the subject of further research.

Switzerland (no data on it)

Table 16 displays the number of domains attacked in each year along with their associated countries present in the final dataset. Our dataset spans from 2009 until early 2013. Secondary data for measuring the independent variables is only available until 2012, so in our bivariate analysis, 2013 is not covered.

TABLE 16- ATTACKED DOMAINS PER REGION PER YEAR

Year	All Domains	All Countries	EU Domains	EU Countries
2009	807	45	346	23
2010	918	51	387	23
2011	869	49	376	23
2012	816	51	359	23

7.5.5 POOLED DATA ON ALL YEARS, VERSUS FOCUSING ON A SINGLE YEAR

A similar argument can be made regarding time: it could happen that a pattern that holds across all years might not hold for a particular year, say 2012, and vice versa. Therefore, trying different combinations of years rather than just analyzing the dataset for all years might provide us with more information. For this reason, we use two versions of the datasets in the analysis: one pooled, containing all years, and another with only 2012 data.

7.6 TEST OF INDIVIDUAL HYPOTHESIS

This section consists of sub-sections where each hypothesis is tested. In the first sub-section normality of each of the variables will be investigated to determine which kind of test should be performed using different variables (parametric/non-parametric tests). At the end of each sub-section under heading 'finding', the results for the four different sets are compared and reported. The significance level of the analysis is set at 0.05. A visual presentation is provided only for tests that have a significant level lower than this.

7.6.1 TEST OF NORMALITY

In this section we are going to test which of the variables follows normal distribution to determine what statistical test to perform further for analysis. The variables consist of a dependent variable (domain and country level) and three non-categorical independent variables. The hypothesis that is going to be tested is as follow:

The sample population is normally distributed.

Dependent Variable

Domain Level

Result: H0 is rejected and the variable 'attacks' is not normally distributed.

Tests of Normality

	Kolmogorov-Smirnov ^a			Shapiro-Wilk		
	Statistic	df	Sig.	Statistic	df	Sig.
attacks	.312	3405	.000	.534	3405	.000

a. Lilliefors Significance Correction

FIGURE 46- RESULT OF TEST OF NORMALITY

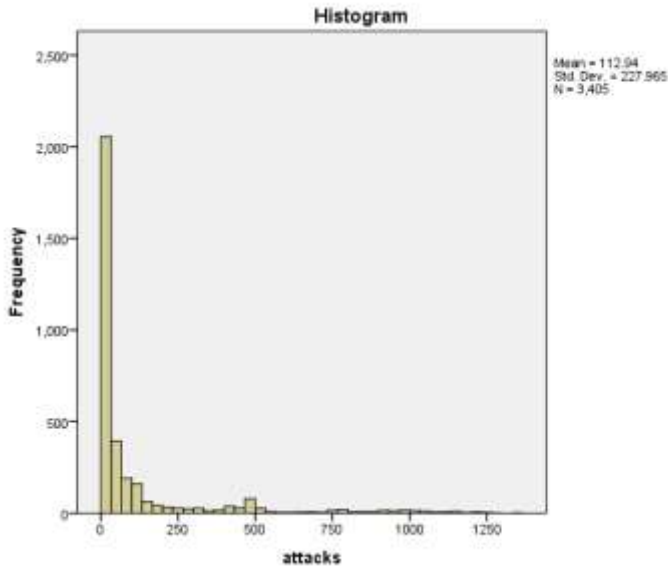


FIGURE 47- HISTOGRAM OF VARIABLE 'ATTACKS'

Descriptives

			Statistic	Std. Error
attacks	Mean		112.94	3.907
	95% Confidence Interval for Mean	Lower Bound	105.28	
		Upper Bound	120.60	
	5% Trimmed Mean		73.00	
	Median		20.00	
	Variance		51967.843	
	Std. Deviation		227.965	
	Minimum		1	
	Maximum		1359	
	Range		1358	
	Interquartile Range		75	
	Skewness		2.913	.042
	Kurtosis		8.378	.084

FIGURE 48- DESCRIPTIVES OF VARIABLE 'ATTACKS'

Country Level

Result: H0 is rejected and the variable 'sum of attacks' is not normally distributed.

Tests of Normality

	Kolmogorov-Smirnov ^a			Shapiro-Wilk		
	Statistic	df	Sig.	Statistic	df	Sig.
sum_attacks	.391	191	.000	.369	191	.000

a. Lilliefors Significance Correction

FIGURE 49- RESULT OF KS TEST OF NORMALITY

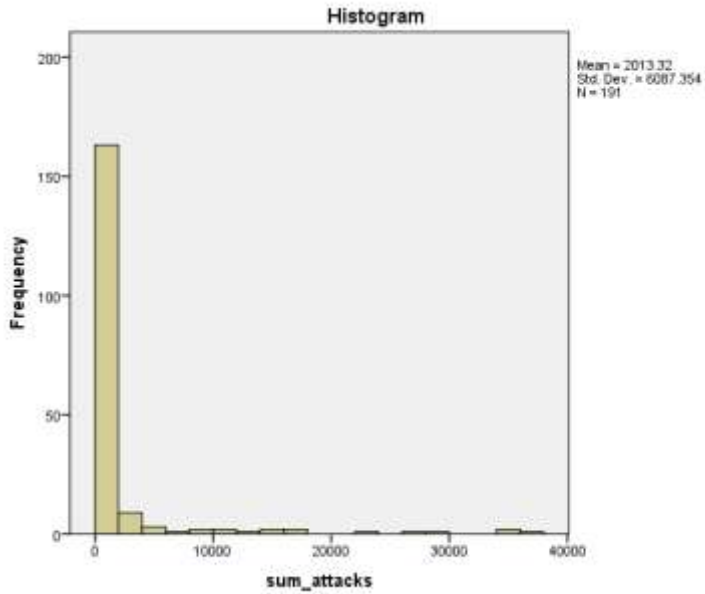


FIGURE 50- HISTOGRAM OF VARIABLE 'SUM OF ATTACKS'

Descriptives

		Statistic	Std. Error
sum_attacks	Mean	2013.32	440.465
	95% Confidence Interval for Mean	Lower Bound 1144.49	
		Upper Bound 2882.16	
	5% Trimmed Mean	812.14	
	Median	63.00	
	Variance	37055876.86	
	Std. Deviation	6087.354	
	Minimum	1	
	Maximum	37309	
	Range	37308	
	Interquartile Range	523	
	Skewness	4.136	.176
	Kurtosis	17.821	.350

FIGURE 51- DESCRIPTIVES OF VARIABLE 'SUM OF ATTACKS'

Independent Variables

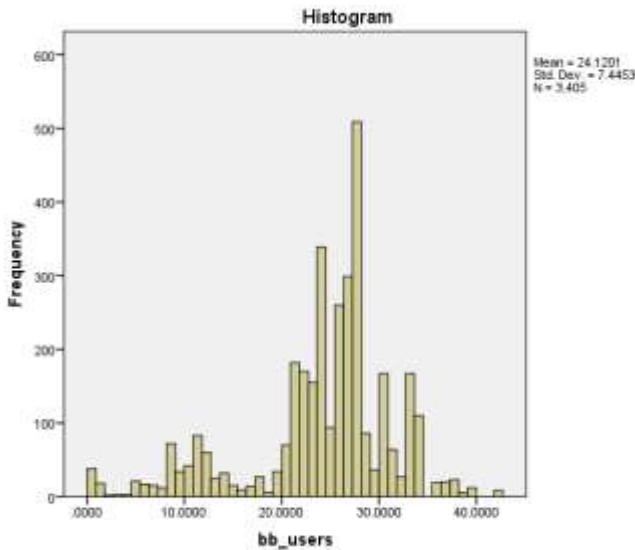
Broadband users

Result: H0 is rejected and the variable 'broadband users' is not normally distributed.

Tests of Normality

	Kolmogorov-Smirnov ^a			Shapiro-Wilk		
	Statistic	df	Sig.	Statistic	df	Sig.
bb_users	.154	3405	.000	.916	3405	.000

a. Lilliefors Significance Correction



Descriptives

			Statistic	Std. Error
bb_users	Mean		24.120138	.1275915
	95% Confidence Interval for Mean	Lower Bound	23.869974	
		Upper Bound	24.370302	
	5% Trimmed Mean		24.528454	
	Median		25.460000	
	Variance		55.432	
	Std. Deviation		7.4452653	
	Minimum		.0100	
	Maximum		41.8600	
	Range		41.8500	
	Interquartile Range		5.9300	
	Skewness		-1.021	.042
	Kurtosis		1.116	.084

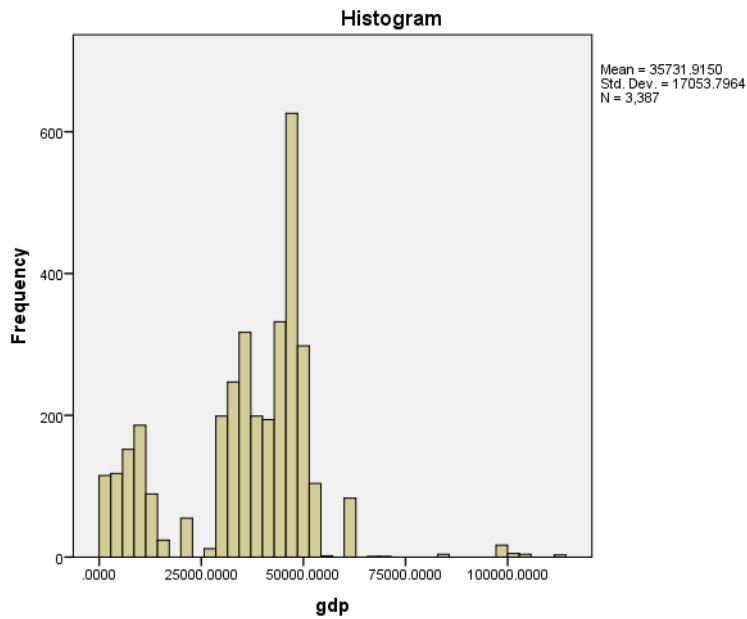
GDP

Result: H0 is rejected and the variable 'GDP' is not normally distributed.

Tests of Normality

	Kolmogorov-Smirnov ^a			Shapiro-Wilk		
	Statistic	df	Sig.	Statistic	df	Sig.
gdp	.132	3387	.000	.889	3387	.000

a. Lilliefors Significance Correction



Descriptives

			Statistic	Std. Error
gdp	Mean		35731.91498	293.0309382
	95% Confidence Interval for Mean	Lower Bound	35157.37952	
		Upper Bound	36306.45043	
	5% Trimmed Mean		35907.81902	
	Median		39659.05886	
	Variance		2.908E8	
	Std. Deviation		17053.79640	
	Minimum		129.3649	
	Maximum		114231.7508	
	Range		114102.3859	
	Interquartile Range		17184.0663	
	Skewness		-.221	.042
	Kurtosis		1.062	.084

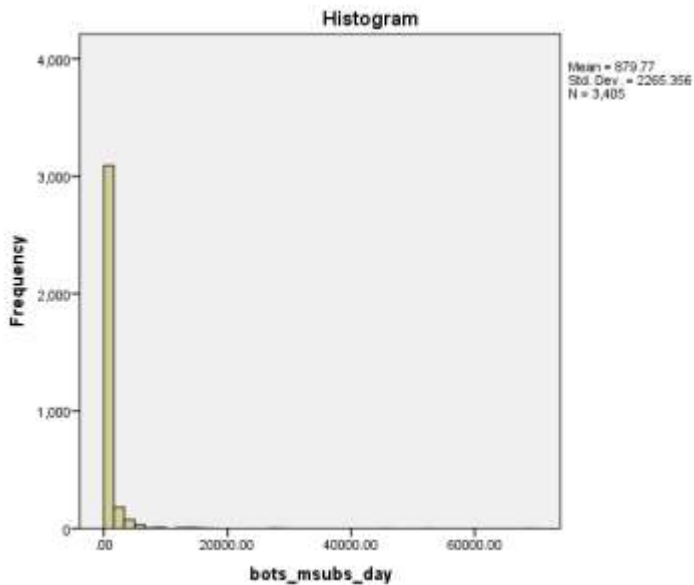
Infection Rate

Result: H0 is rejected and the variable 'Infection Rate' is not normally distributed.

Tests of Normality

	Kolmogorov-Smirnov ^a			Shapiro-Wilk		
	Statistic	df	Sig.	Statistic	df	Sig.
bots_msubs_day	.353	3405	.000	.261	3405	.000

a. Lilliefors Significance Correction



Descriptives

			Statistic	Std. Error
bots_msubs_day	Mean		879.7712	38.82200
	95% Confidence Interval for Mean	Lower Bound	803.6544	
		Upper Bound	955.8880	
	5% Trimmed Mean		604.5452	
	Median		448.0000	
	Variance		5131837.924	
	Std. Deviation		2265.35603	
	Minimum		23.00	
	Maximum		68405.00	
	Range		68382.00	
	Interquartile Range		951.00	
	Skewness		16.038	.042
	Kurtosis		373.883	.084

Finding

As the result of normality tests performed above indicates, none of the variables follow normal distribution. Therefore, from now on in the report, non-parametric tests will be performed on the variables displays above.

7.6.2 COMPARISON OF MEANS

Hypothesis 1

Banks with word 'bank' in their domain names are more popular among cybercriminals.

This hypothesis is tested by comparing the median of the group that covers domain names with word 'bank' and the group that contains domains without word 'bank'.

1. *Year=2012, Region= EU*

Result: Ho cannot be rejected and the medians are more or less the same in both groups.

	with_bank	N	Mean Rank	Sum of Ranks
attacks	0	279	185.01	51618.00
	1	80	162.53	13002.00
	Total	359		

	attacks
Mann-Whitney U	9762.000
Wilcoxon W	13002.000
Z	-1.711
Asymp. Sig. (2-tailed)	.087

a. Grouping Variable:
with_bank

FIGURE 52- RESULT OF THE COMPARISON OF MEANS TEST FOR H1

2. *Year=2012, Region= all regions*

Result: Ho cannot be rejected and the medians are more or less the same in both groups.

	with_bank	N	Mean Rank	Sum of Ranks
attacks	0	633	411.99	260789.50
	1	182	394.12	71730.50
	Total	815		

	attacks
Mann-Whitney U	55077.500
Wilcoxon W	71730.500
Z	-.905
Asymp. Sig. (2-tailed)	.365

a. Grouping Variable:
with_bank

FIGURE 53- RESULT OF THE COMPARISON OF MEANS TEST FOR H1

3. *Year=all years, Region= EU*

Result: Ho is rejected and the medians are not the same and domains with 'bank' word are less popular among cybercriminals (reverse of expected direction).

Ranks				Test Statistics ^a		
	with_bank	N	Mean Rank	Sum of Ranks	attacks	
attacks	0	1184	752.72	891223.00	Mann-Whitney U	146553.000
	1	284	658.53	187023.00	Wilcoxon W	187023.000
Total		1468			Z	-3.364
					Asymp. Sig. (2-tailed)	.001

a. Grouping Variable: with_bank

FIGURE 54- RESULT OF THE COMPARISON OF MEANS TEST FOR H1

4. *Year=all years, Region= all regions*

Result: Ho cannot be rejected and the medians are more or less the same in both groups.

Ranks				Test Statistics ^a		
	with_bank	N	Mean Rank	Sum of Ranks	attacks	
attacks	0	2673	1704.70	4556659.00	Mann-Whitney U	973778.000
	1	732	1696.80	1242056.00	Wilcoxon W	1242056.000
Total		3405			Z	-.193
					Asymp. Sig. (2-tailed)	.847

a. Grouping Variable: with_bank

FIGURE 55- RESULT OF THE COMPARISON OF MEANS TEST FOR H1

Finding

Based on the result of our analysis, the median of attacks are more or less the same in domains with 'bank' word in their domain name and those without word 'bank' in their domain names. Thus, considering the fact that the idea for this variable has primarily been discussed by the security experts working on the same field, we can conclude that attackers actually act more professionally. That is, they have a good understating of market they are working in or have access to people with knowledge about popular local domains.

Hypothesis 2

A bank is more popular among cybercriminals as a target when it offers online banking in English.

1. *Year=2012, Region= EU*

Result: H0 is rejected and domains with English webpages are more popular among cybercriminals.

Ranks					Test Statistics ^a	
	with_en	N	Mean Rank	Sum of Ranks		attacks
attacks	0	187	163.76	30623.50	Mann-Whitney U	13045.500
	1	172	197.65	33996.50	Wilcoxon W	30623.500
	Total	359			Z	-3.096
					Asymp. Sig. (2-tailed)	.002

a. Grouping Variable: with_en

FIGURE 56- RESULT OF THE COMPARISON OF MEANS TEST FOR H2

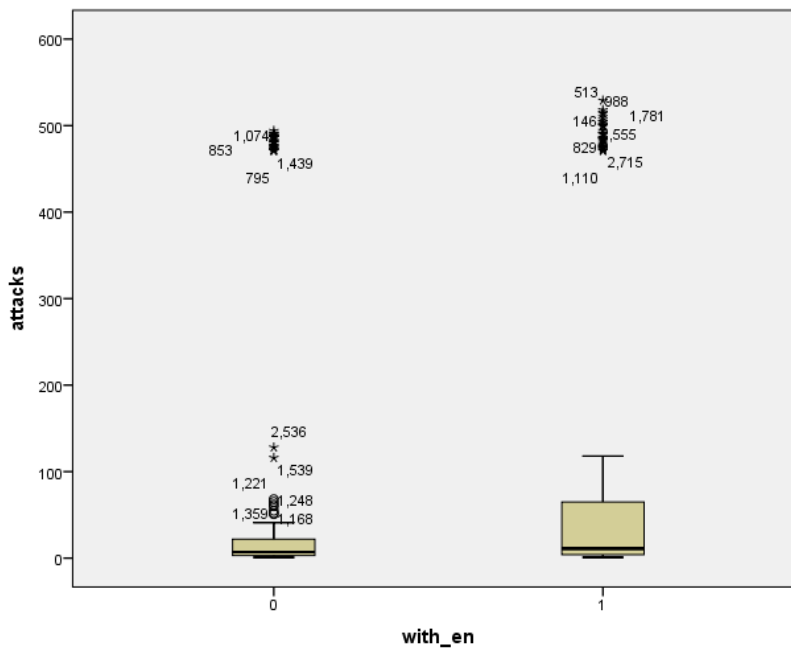


FIGURE 57- BOX PLOT OF SOM OF ATTACKS WITHIN TWO GROUPS OF WITH/WITHOUT ENGLISH WEBPAGES

2. *Year=2012, Region= all regions*

Result: H0 cannot be rejected and domains with English webpages are as popular among cybercriminals as domains without English webpages.

Ranks				Test Statistics ^a		
	with_en	N	Mean Rank	Sum of Ranks	attacks	
attacks	0	229	396.82	90872.50	Mann-Whitney U	64537.500
	1	586	412.37	241647.50	Wilcoxon W	90872.500
Total				815	Z	-.850
					Asymp. Sig. (2-tailed)	.395

a. Grouping Variable: with_en

FIGURE 58- RESULT OF THE COMPARISON OF MEANS TEST FOR H2

3. *Year=all years, Region= EU*

Result: H0 is rejected and domains with English webpages are more popular among cybercriminals than domains without English webpages.

Ranks				Test Statistics ^a		
	with_en	N	Mean Rank	Sum of Ranks	attacks	
attacks	0	704	692.84	487757.00	Mann-Whitney U	239597.000
	1	764	772.89	590489.00	Wilcoxon W	487757.000
Total				1468	Z	-3.616
					Asymp. Sig. (2-tailed)	.000

a. Grouping Variable: with_en

FIGURE 59- RESULT OF THE COMPARISON OF MEANS TEST FOR H2

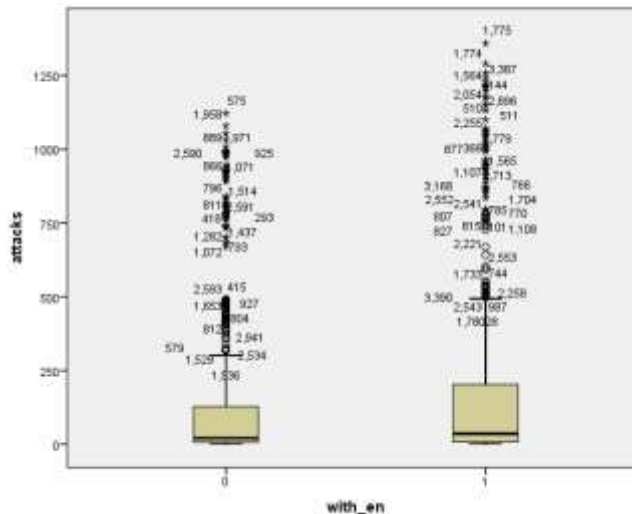


FIGURE 60- BOX PLOT OF SOM OF ATTACKS WITHIN TWO GROUPS OF WITH/WITHOUT ENGLISH WEBPAGES

4. *Year=all years, Region= all regions*

Result: H0 cannot be rejected and domains with English webpages are as popular among cybercriminals as domains without English webpages.

Ranks				Test Statistics ^a		
	with_en	N	Mean Rank	Sum of Ranks	attacks	
attacks	0	939	1712.12	1607676.50	Mann-Whitney U	1149227.500
	1	2466	1699.53	4191038.50	Wilcoxon W	4191038.500
	Total	3405			Z	-.334
					Asymp. Sig. (2-tailed)	.738

a. Grouping Variable: with_en

FIGURE 61- RESULT OF THE COMPARISON OF MEANS TEST FOR H2

Finding

Based on the result of our analysis, we can recognize the average of attacks on domains with English webpages is more than those without, in the EU countries. However, this is not the case when looking at all countries. The reason might be explained by the fact that the EU region contains countries with a wide variety of languages, which have similarities in terms of infrastructure, etc. So if we consider that everything rather than language is more or less the same for attacking banks in EU region, then it is logical that attackers prefer domains with English options over other domains. This is simply because in this way they can attack a wide variety of domains in different countries in EU with less effort (in term of language knowledge).

7.6.3 MEASURES OF ASSOCIATION

Hypothesis 3

Banks in countries with higher broadband penetration are more popular targets among cybercriminals.

1. *Year=2012, Region= EU*

Result: H0 cannot be rejected and domains in countries with higher rate of bb-penetration are attacked as much as other domains.

Correlations			sum_attacks	bb_users
Spearman's rho	sum_attacks	Correlation Coefficient	1.000	.404
		Sig. (2-tailed)	.	.069
		N	21	21
	bb_users	Correlation Coefficient	.404	1.000
		Sig. (2-tailed)	.069	.
		N	21	21

FIGURE 62- RESULT OF SPEARMAN'S BIVARIATE CORRELATION TEST FOR H3

2. *Year=2012, Region= all regions*

Result: H0 is rejected and domains in countries with higher rate of broadband penetration are attacked more.

Correlations

			sum_attacks	bb_users
Spearman's rho	sum_attacks	Correlation Coefficient	1.000	.467**
		Sig. (2-tailed)	.	.001
		N	50	50
	bb_users	Correlation Coefficient	.467**	1.000
		Sig. (2-tailed)	.001	.
		N	50	50

** . Correlation is significant at the 0.01 level (2-tailed).

FIGURE 63- RESULT OF SPEARMAN'S BIVARIATE CORRELATION TEST FOR H3

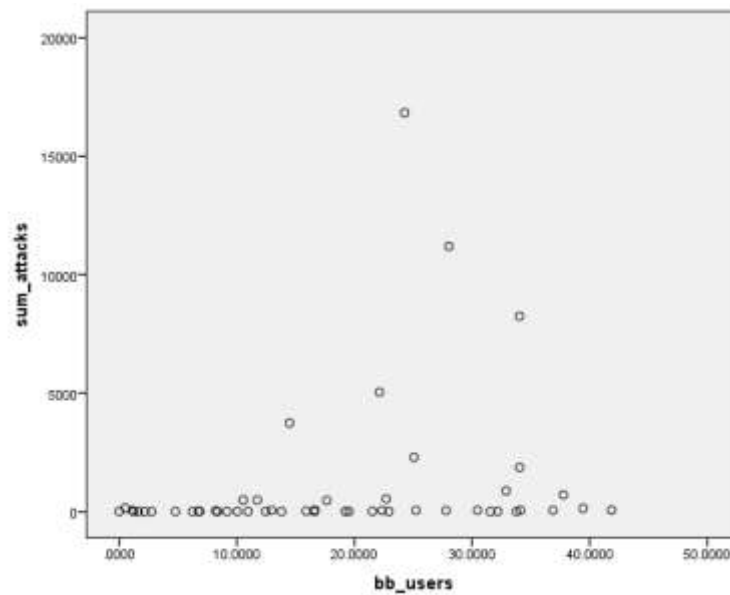


FIGURE 64- SCATTER PLOT OF BROADBAND USERS VERSUS SUM OF ATTACKS ON DIFFERENT COUNTRIES

3. Year=all years, Region= EU

Result: H0 cannot be rejected and domains in countries with higher rate of bb-penetration are attacked as much as other domains.

Correlations

			sum_attacks	bb_users
Spearman's rho	sum_attacks	Correlation Coefficient	1.000	.161
		Sig. (2-tailed)	.	.154
		N	80	80
	bb_users	Correlation Coefficient	.161	1.000
		Sig. (2-tailed)	.154	.
		N	80	80

FIGURE 65- RESULT OF SPEARMAN'S BIVARIATE CORRELATION TEST FOR H3

4. Year=all years, Region= all regions

Result: H0 is rejected and domains in countries with higher rate of bb-penetration are attacked more.

Correlations

			sum_attacks	bb_users
Spearman's rho	sum_attacks	Correlation Coefficient	1.000	.281**
		Sig. (2-tailed)	.	.000
		N	191	191
	bb_users	Correlation Coefficient	.281**	1.000
		Sig. (2-tailed)	.000	.
		N	191	191

** . Correlation is significant at the 0.01 level (2-tailed).

FIGURE 66- RESULT OF SPEARMAN'S BIVARIATE CORRELATION TEST FOR H3

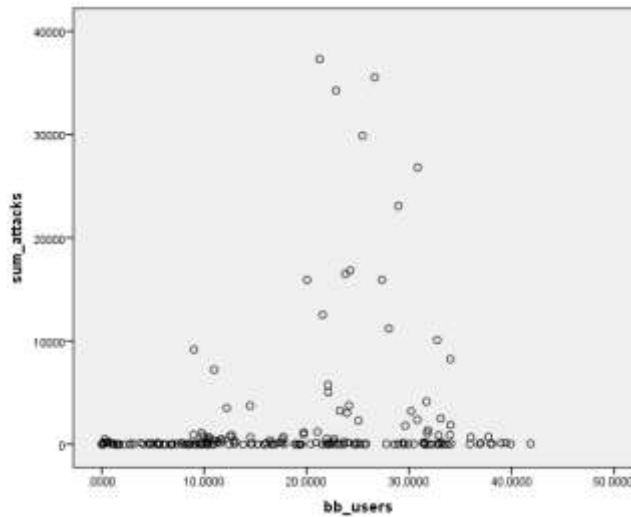


FIGURE 67- SCATTER PLOT OF BROADBAND USERS VERSUS SUM OF ATTACKS ON DIFFERENT COUNTRIES

Finding

Based on the results presented in the above four combinations of data, It can be concluded that as the number of broadband users increase as, they become more valuable targets because they can be potential targets looking from cybercriminals' perspective. However, the statement holds true mostly for the global region.

Hypothesis 4

Banks in countries with higher GDP are more popular attack targets among cybercriminals.

1. *Year=2012, Region= EU*

Result: H0 cannot be rejected and domains in countries with higher GDP are attacked as much as countries with lower GDP.

Correlations

			sum_attacks	gdp
Spearman's rho	sum_attacks	Correlation Coefficient	1.000	.228
		Sig. (2-tailed)	.	.320
		N	21	21
	gdp	Correlation Coefficient	.228	1.000
		Sig. (2-tailed)	.320	.
		N	21	21

FIGURE 68- RESULT OF SPEARMAN'S BIVARIATE CORRELATION TEST FOR H4

2. *Year=2012, Region= all regions*

Result: H0 is rejected and domains in countries with higher GDP are attacked more.

Correlations

			sum_attacks	gdp
Spearman's rho	sum_attacks	Correlation Coefficient	1.000	.407**
		Sig. (2-tailed)	.	.006
		N	50	44
	gdp	Correlation Coefficient	.407**	1.000
		Sig. (2-tailed)	.006	.
		N	44	44

** . Correlation is significant at the 0.01 level (2-tailed).

FIGURE 69- RESULT OF SPEARMAN'S BIVARIATE CORRELATION TEST FOR H4

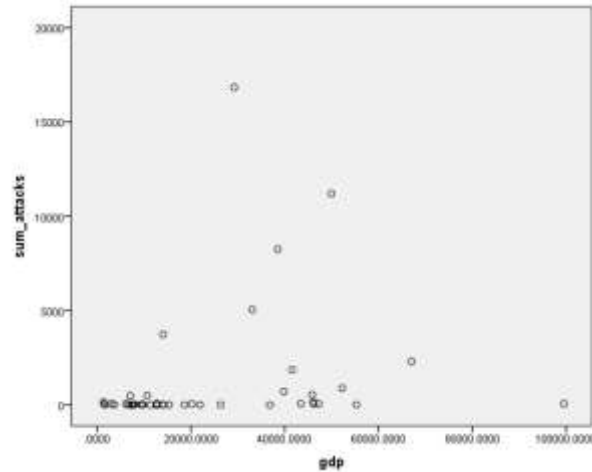


FIGURE 70-SCATTER PLOT OF GDP VERSUS SUM OF ATTACKS ON DIFFERENT COUNTRIES

3. *Year=all years, Region= EU*

Result: H0 cannot be rejected and domains in countries with higher GDP are attacked as much as countries with lower GDP.

Correlations

			sum_attacks	gdp
Spearman's rho	sum_attacks	Correlation Coefficient	1.000	.135
		Sig. (2-tailed)	.	.231
		N	80	80
	gdp	Correlation Coefficient	.135	1.000
		Sig. (2-tailed)	.231	.
		N	80	80

FIGURE 71-RESULT OF SPEARMAN'S BIVARIATE CORRELATION TEST FOR H4

4. *Year=all years, Region= all regions*

Result: H0 is rejected and domains in countries with higher GDP are attacked more.

Correlations

			sum_attacks	gdp
Spearman's rho	sum_attacks	Correlation Coefficient	1.000	.270**
		Sig. (2-tailed)	.	.000
		N	191	183
	gdp	Correlation Coefficient	.270**	1.000
		Sig. (2-tailed)	.000	.
		N	183	183

** . Correlation is significant at the 0.01 level (2-tailed).

FIGURE 72-RESULT OF SPEARMAN'S BIVARIATE CORRELATION TEST FOR H4

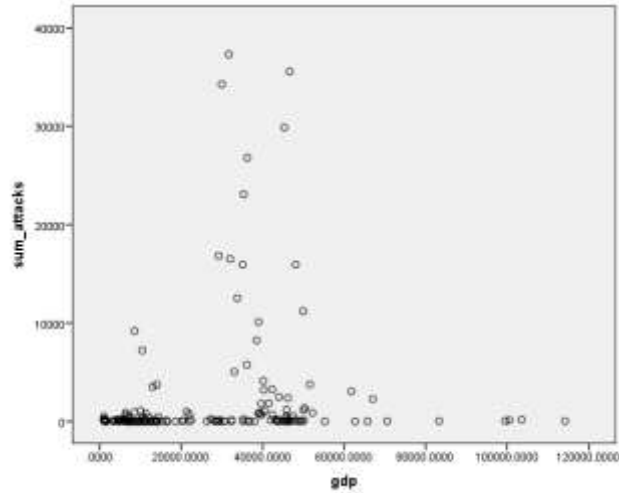


FIGURE 73- SCATER PLOT OF GDP VERSUS SUM OF ATTACKS ON DIFFERENT COUNTRIES

Finding

We expected cybercriminals to have an intuition about the statistics in regards to some country related statistics based on our conceptual model. For instance they more or less know which countries have high GDP. However, as it is obvious from the above tests, the hypothesis holds true globally but not when comparing just EU countries. The reason boils down to the sample variation; that is GDP matters, but when they are within a specific range, than other factors are more important. The reason is explained in more details in 'EU vs. Global' sub-section.

Hypothesis 5

Banks in countries with higher rate of infection are more popular attack targets among cybercriminals.

1. *Year=2012, Region= EU*

Result: H0 cannot be rejected and domains in countries with higher infection rate are attacked as countries with lower infection rate.

Correlations

			sum_attacks	bots_msubs_day
Spearman's rho	sum_attacks	Correlation Coefficient	1.000	-.143
		Sig. (2-tailed)	.	.536
		N	21	21
	bots_msubs_day	Correlation Coefficient	-.143	1.000
		Sig. (2-tailed)	.536	.
		N	21	21

FIGURE 74- RESULT OF SPEARMAN'S BIVARIATE CORRELATION TEST FOR H5

2. *Year=2012, Region= all regions*

Result: H0 is rejected and domains in countries with higher infection rate are attacked more than others (Reverse of the expected direction).

Correlations

			sum_attacks	bots_msubs_day
Spearman's rho	sum_attacks	Correlation Coefficient	1.000	-.396**
		Sig. (2-tailed)	.	.004
		N	50	50
	bots_msubs_day	Correlation Coefficient	-.396**	1.000
		Sig. (2-tailed)	.004	.
		N	50	50

** . Correlation is significant at the 0.01 level (2-tailed).

FIGURE 75- RESULT OF SPEARMAN'S BIVARIATE CORRELATION TEST FOR H5

3. *Year=all years, Region= EU*

Result: H0 cannot be rejected and domains in countries with higher infection rate are attacked as countries with lower infection rate.

Correlations			sum_attacks	bots_msubs_day
Spearman's rho	sum_attacks	Correlation Coefficient	1.000	.122
		Sig. (2-tailed)	.	.280
		N	80	80
	bots_msubs_day	Correlation Coefficient	.122	1.000
		Sig. (2-tailed)	.280	.
		N	80	80

FIGURE 76- RESULT OF SPEARMAN'S BIVARIATE CORRELATION TEST FOR H5

4. *Year=all years, Region= all regions*

Result: H0 cannot be rejected and domains in countries with higher infection rate are attacked as countries with lower infection rate.

Correlations			sum_attacks	bots_msubs_day
Spearman's rho	sum_attacks	Correlation Coefficient	1.000	-.079
		Sig. (2-tailed)	.	.275
		N	191	191
	bots_msubs_day	Correlation Coefficient	-.079	1.000
		Sig. (2-tailed)	.275	.
		N	191	191

FIGURE 77- RESULT OF SPEARMAN'S BIVARIATE CORRELATION TEST FOR H5

Finding

We expected that cybercriminals have an intuition about the statistics in regards to the infection rate of different countries. For instance they more or less know which countries have high infection rate. However, the result of our analysis indicates that criminals may not make their attack decisions completely based on such statistics but it is more probable that they use any of their chance for executing an online attack on a wide range of target domains. The results confirm our conclusion as none of them has correlation with number of attacks.

7.6.4 CORRELATION MATRIX

Presented below is the correlation matrix between all the independent variables.

Variables: with bank – with English

Correlations

			with_bank	with_en
Spearman's rho	with_bank	Correlation Coefficient	1.000	.021
		Sig. (2-tailed)	.	.230
		N	3405	3405
	with_en	Correlation Coefficient	.021	1.000
		Sig. (2-tailed)	.230	.
		N	3405	3405

FIGURE 78- RESULT OF SPEARMAN'S CORRELATION TEST FOR CATEGORICAL INDEPENDENT VARIABLES

Variables: Broad Band Penetration- GDP- Infection Rate

Correlations

			bb_users	gdp	bots_msubs_ day
Spearman's rho	bb_users	Correlation Coefficient	1.000	.791**	-.640**
		Sig. (2-tailed)	.	.000	.000
		N	191	183	191
	gdp	Correlation Coefficient	.791**	1.000	-.480**
		Sig. (2-tailed)	.000	.	.000
		N	183	183	183
	bots_msubs_ day	Correlation Coefficient	-.640**	-.480**	1.000
		Sig. (2-tailed)	.000	.000	.
		N	191	183	191

** . Correlation is significant at the 0.01 level (2-tailed).

FIGURE 79- RESULT OF SPEARMAN'S BIVARIATE CORRELATION TEST FOR NON-GEORAPHICAL INDEPENDENT VARIABLES

Summary

Table 17 displays the summary of the findings in this section.

TABLE 17- SUMMARY OF THE FINDINGS

Hypothesis	Statistical test	2012,EU	2012,all regions	All years, EU	Pooled data	Level of Analysis	Verdict EU set	Verdict global set
1 Banks with 'bank' word	Comparison of means	Sig=0.087 M No>M Yes	Sig=0.365 M No>M Yes	Sig=0.001 M No>M Yes	Sig=0.847 M No>M Yes	Domain	Rejected	Rejected
2 Banks with English webpage	Comparison of means	Sig=0.002** M Yes>M No	Sig=0.395 M Yes>M No	Sig=0.000*** M Yes>M No	Sig=0.738 M No>M Yes	Domain	Accepted	Rejected
3 Effects of BB penetration	Spearman's	Sig=0.069 Corr=0.404	Sig=0.001** Corr=0.467	Sig=0.151 Corr=0.161	Sig=0.000*** Corr=0.281	Country	Rejected	Accepted
4 Effects of GDP	Spearman's	Sig=0.320 Corr=0.228	Sig=0.006** Corr=0.407	Sig=0.135 Corr=0.231	Sig=0.000*** Corr=0.270	Country	Rejected	Accepted
5 Effects of infection rate	Spearman's	Sig=0.536 Corr= -0.143	Sig=0.004** Corr= - 0.396	Sig=0.280 Corr=0.122	Sig=0.275 Corr= - 0.079	Country	Rejected	Rejected

7.6.5 EU VERSUS GLOBAL

In this section we will review our findings about hypothesis 2, 3 and 4 and will elaborate more on them. As Table 17 displays, the results of test performed on these hypothesis holds true only for specific regions.

Starting from hypothesis 2, we saw that domains with English domain pages within the EU region have been attacked more on average comparing to the global region. We explained that, the underlying reason of the above result might be explained by the fact that the EU region contains countries with a wide variety of languages, which have similarities in terms of infrastructure, etc. This means that if we suppose attackers require more or less the same infrastructures for attacking EU region banks, apart from language, then it is logical if they attack domains with English language option more than other domains because it requires less effort. However, the hypothesis does not hold true globally because the globally countries are the variety of infrastructures in different countries around the world increase and thus the pattern cannot be seen anymore.

The same arguments can be discussed for hypothesis 3 and 4. Figure 80 displays the result of descriptive statistics for two variables 'broadband users' and 'GDP' within the EU region. Figure 81 displays result of descriptive statistics for two variables 'broadband users' and 'GDP' globally. Comparing the two figures with each other using 'coefficient of variance' or CV, we will realize that variety is more in both of the variables within the global region. Likewise, considering the fact that 'GDP' and 'broadband users' are infrastructure variables, it is logical that attackers take these variables into account, or these variables become more important where variety is more comparing to the EU region where variety is less. This finding can actually be marked as one of the most important findings in our bivariate analysis.

EU Region

CV (bb users) = 0.3

CV (GDP) = 0.62

Descriptive Statistics

	N	Minimum	Maximum	Mean	Std. Deviation
bb_users	80	12.5400	39.4400	24.604125	7.4561596
gdp	80	6334.6821	114231.7508	34351.32831	21545.46450
Valid N (listwise)	80				

FIGURE 80- DESCRIPTIVE STATISTICS WITHIN THE EU REGION

Global Region

CV (bb users) = 0.62

CV (GDP) = 0.83

Descriptive Statistics

	N	Minimum	Maximum	Mean	Std. Deviation
bb_users	191	.0100	41.8600	18.076754	11.2266051
gdp	183	949.1166	114231.7508	26062.69456	21708.11833
Valid N (listwise)	183				

FIGURE 81- DESCRIPTIVE STATISTICS FOR THE GLOBAL REGION

7.7 DECISION TREE ANALYSIS

Overview

Decision tree analysis is a data mining technique defined as the process of discovering patterns in data (Witten & Frank, 2005). Due to the explorative nature of our research, this method might provide additional insights. Additional insight in the sense that we can understand which of the independent variables has more influence on our dependent variables recursively. Moreover, this method can be used to validate the result of our statistical measurements as well (those variables that are also included in this test). The decision tree test has been performed using the software package WEKA.

7.7.1 BUILDING THE DATASET

Before starting the test, we need to determine what the decision variable is. In other words, what we intend to get as an output from our decision-tree analysis. Given the goal of this research, we set the decision variable to the **'average number of botnets attacking a domain (per week)'**, since we would like to know more about popularity of attack domains amongst cybercriminals.

Executing the decision tree test on WEKA environment requires a specific type of dataset. We have to convert all variables in our dataset into categories. The four variables of *attack category*, *country category (region)*, *with English* and *with bank* will be used in this analysis. We categorize the attacked domains into two groups: rarely attacked and mostly attacked. For the country category variable, the countries in our dataset are divided into the five regions listed in Table 18. The two variables *with English* and *with bank* are divided into Yes/No categories.

TABLE 18- COUNTRIES IN OUR DATASET

Category (Region)	Countries
Western Europe	Germany, United Kingdom, Ireland, Spain, Portugal ,Italy, Finland, France, Switzerland, Belgium, Norway, Luxembourg, Iceland, Sweden, Austria, Netherlands, Greece
Central and Eastern Europe	Russia, Latvia, Turkey, Cyprus, Romania, Poland, Belarus, Czech Republic, Lithuania, Bulgaria, Estonia, Hungary
North America	United States, Canada, Panama, Mexico, Bahamas, Costa Rica
Latin America and Caribbean	Argentina, Chile, Colombia, Venezuela, Peru, Cayman Islands, Antigua and Barbuda, Brazil
Middle East and Africa	Pakistan, United Arab Emirates, Saudi Arabia, South Africa, Azerbaijan, Nigeria, Kuwait
Asia and Pacific	Australia, China, India, Thailand, New Zealand, Singapore, Hong King, Malaysia, Indonesia

7.7.2 RANDOM FOREST TEST

After building the dataset, we are prepared to perform the test. The test that is performed is called **'Random Forest'** classification test and it is set to cross validate the result with 10 folds from data. The Random Forest technique is been chosen since it often performs "best of all" (for "easy-to-understand reasons"). The test results are presented in the box below and in Figure 82 and described thereafter.

Test Summary		
Correctly Classified Instances	701	59.7613 %
Incorrectly Classified Instances	472	40.2387 %
Mean absolute error	0.473	
Root mean squared error	0.491	
Relative absolute error	94.65 %	
Root relative squared error	98.25 %	
Coverage of cases (0.95 level)	99.82 %	
Mean rel. region size (0.95 level)	99.61 %	
Total Number of Instances	1173	

=== Confusion Matrix ===

a b <-- classified as
 360 225 | a = Mostly
 247 341 | b = rarely

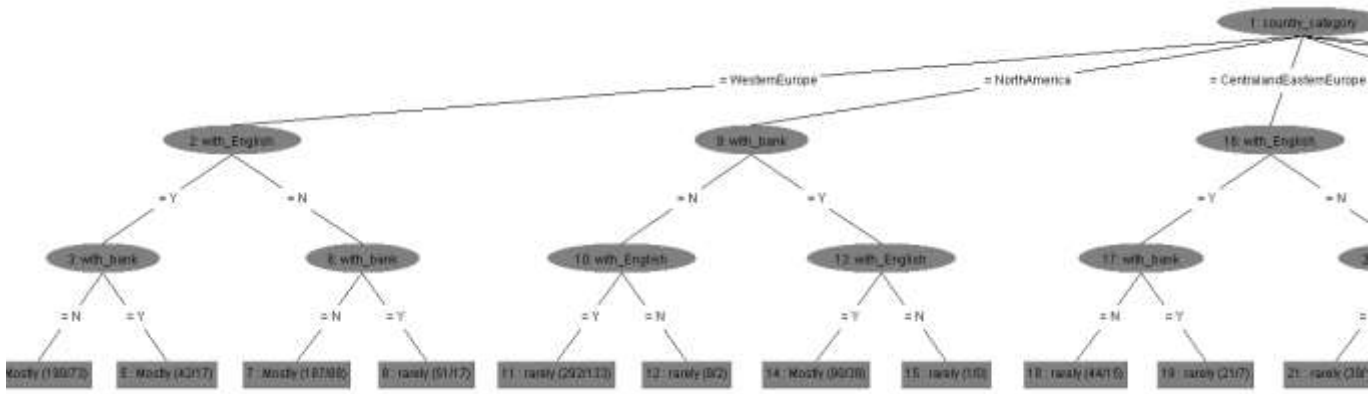


FIGURE 82 – VISUAL TREE (LEFT PART)

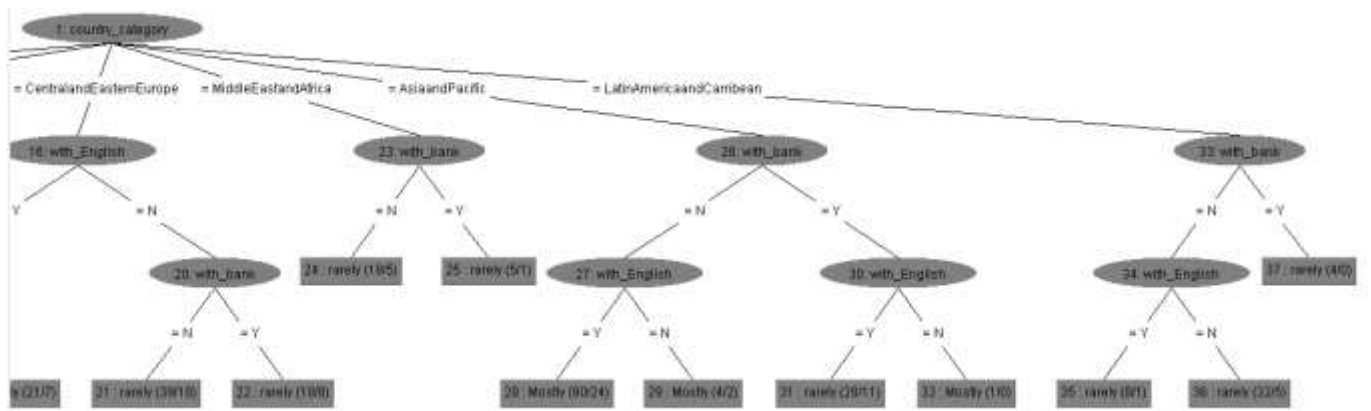


FIGURE 83 – VISUAL TREE (RIGHT PART)

The results indicate:

- 1) Data show little structure/much noise since
 - a. Prediction accuracy is only around 60% (number of correctly classified instances). This is slightly better than throwing coins randomly.
 - b. The decision tree has been fully elaborated: all possible branches are part of it.
- 2) Variable with highest prediction power is 'country_catgeory' or region.
- 3) As it can be mentioned from the visualize three, in certain 'Western Europe' and 'Central and Eastern Europe', the variable 'with English' is the second most important (after country category) and in others the variable 'with_bank'. These findings are in line with the findings of the bivariate analysis: the result of our statistical bivariate analysis indicates that banks with English option in their domain pages are attacked more within the EU region due to large variety of languages available in the EU region for countries with similarities in their infrastructure.

7.8 SUMMARY OF FINDINGS

In this section we will briefly summarize all of the finding of the previous sections in Table 19 :

TABLE 19- SUMMARY OF THE FINDINGS

Test	Result
Descriptive Analysis	<ul style="list-style-type: none">• 7% of domains attacked 90% of times.• A portion of domains (around 100) are consistently attacked by Zeus malware and a much larger group of domain has been under attack less than ten weeks. The fact that some of the targets are quite short-lived might imply that cybercriminals perform a lot of trial & error and experiment with attacking many domains.
Bivariate Analysis	<ul style="list-style-type: none">• Hypothesis 2, 3 and 4 were found to be significant.• Region plays an important role in the results of analysis.• Variable 'with_English' found to be significant in EU. We concluded that it might be case that due to the variation in languages in different countries within the EU region, if we assume that they have more or less same infrastructures, and then attackers would choose to go for banks with English option (requires less attack effort).• Variables 'broadband users' and 'GDP' found to be significant in the global region. We argued that these variables are infrastructure variables and only become important for attackers in determining target of attacks where variety in countries become more. In a region like EU with where countries have more or less same GDP for instance, then GDP do not play an important role in determining the attack target.
Decision Tree	<ul style="list-style-type: none">• The result of decision tree test is slightly better than chance.• The result of decision tree test is in line with result of bivariate analysis where in EU region attribute 'with_English' plays an important role.

Chapter 8 - Discussion and implications

8.1 REVIEWING THE FINDINGS

Since now in the report, we have already done a lot of empirical and theoretical work to answer the research question. In this chapter, we will recall the research question and will answer it by summarizing the empirical and theoretical evidence presented in this report. The main research question was as follows:

Can we extract intelligence on criminal attack patterns and target selection from the files which financial malware use as instructions for their operations?

Our answer to the research question, as well as our contribution in this research is consists of four main parts:

- Linking traditional criminology literature and Routine Activity Theory to online banking crime (extending on the work of Jar (2005).
- Creating a three-dimensional dataset consist of time of attack, target of attack and the attacker using the attack records executed by Zeus financial malware.
- Systematically enumerating through options for extracting intelligence from Zeus malware data, including looking at the merits and pitfalls of each.
- Conducting empirical analysis (statistical analysis, decision tree analysis) on the extracted data as a proof of concept.

In the next section, we will review the methodological and empirical findings, then present a summary, and provide recommendations. The recommendations will without doubt be subject to debate, and we look at some of those in 'discussion' section. The chapter finishes by discussing the limitations of the research and looking at future research.

8.1.1 A METHOD FOR EXPRESSING POPULARITY OF TARGETS

As discussed in chapter five, we argued that simply counting the number of times a domain name is seen in the configuration files cannot be a good indication of the popularity of target domains among cybercriminals, as it may over/under report the actual counts. That is because the number of times a domain is seen in a configuration file does not necessarily mean that the domain has been under attack in that specific period. It could be that , a configuration file get updated, due to different reasons such as existence of possible mistakes in the configuration file, including new targets, change of the botnet's password etc.; whereas, some cybercriminals may not update their configuration files often because the attack might continue to work without any change being necessary. Thus, in this way, if we simply count the number of times a domain is seen in the configuration files, domains attacked by criminal groups that update their configuration files more often would get higher ranks of popularity while in practice that may not be the case.

Therefore, in this research we tried to find a new metric to express raw counts in terms of something that would be able to articulate relative popularity of target among cybercriminals from Zeus configuration files. For this, we looked at 10 intermediate variables being extracted based on Routine Activity Theory (RAT) three categories. After studying the variables individually, we suggest a new metric for determining the extent of a domain's popularity as a target among cybercriminals: calculating the '**(average) number of botnets attacking a domain (per week)**', over a year.

$$\text{Number of attacks on a domain over a year} = \sum_{k=1}^{52} \text{botnets mentioning domain in their config file in week (k)}$$

With the variable we are suggesting, we actually are counting only one configuration file per week for each botnet. In this way, we are eliminating the limitations that malware configuration updates may cause. How different the previous method and current method determine the number of attacks performed by Zeus malware on domains can be understood by comparing the two graphs below (Figure 84 and Figure 85).

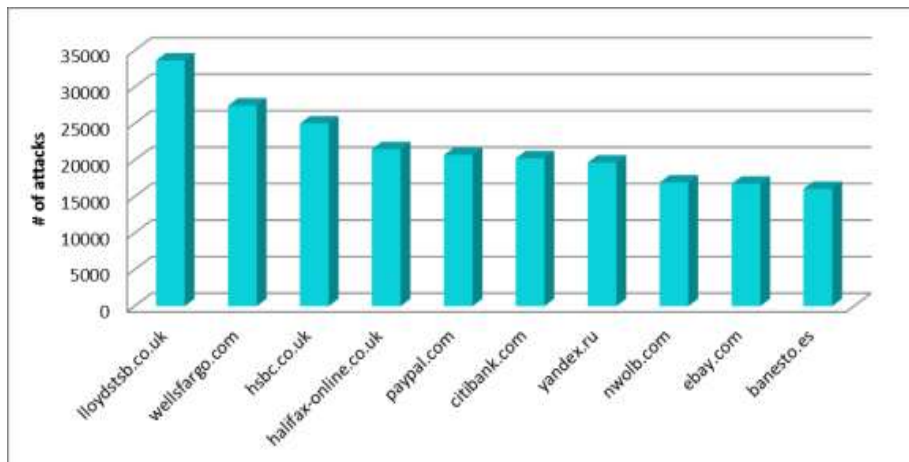


FIGURE 84- TOP-10 ATTACKED DOMAINS (OLD METRIC)

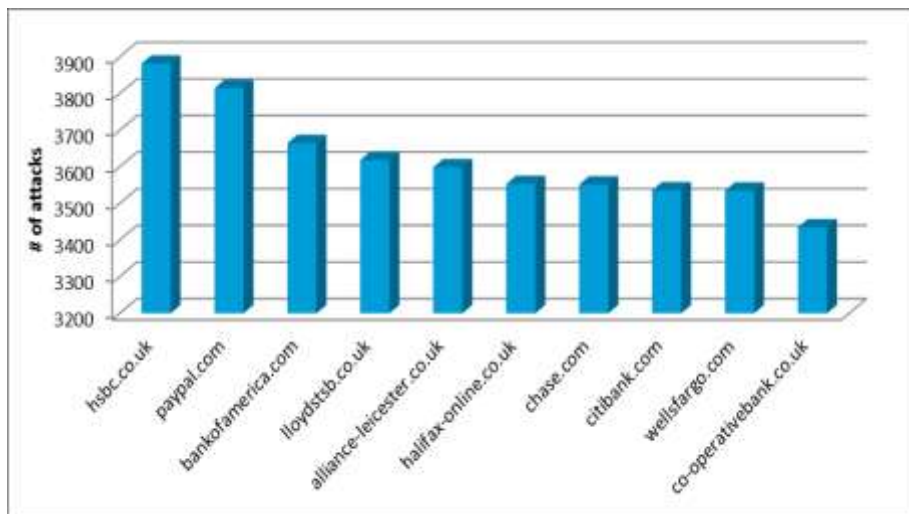


FIGURE 85- TOP-10 ATTACKED DOMAINS (NEW METRIC)

8.1.2 INTELLIGENCE EXTRACTED FROM ZEUS MALWARE CONFIGURATION FILES

We looked at 10643 out of 77899 configuration files of Zeus financial malware. From the metadata of the configuration files, we extract *infection time* of each attack and *login key* as an indication of each botnet. From the web inject section of the configuration files; we extract *domain names* from the injected targeted URLs. After analyzing the preliminary information being extracted from the files via the framework built upon RAT, our findings are as follows:

- ❖ Our dataset confirms the 80-20 power law distribution: a smaller group of domains (15%) attract almost 90% of the attacks.
- ❖ There are a small number of domains (88 domains- **'always-attacked'** group) that are attacked over the whole time span of our data from Jan 2009 to March 2013.
- ❖ There is a large group of domains (1170 domains- **'rarely-attacked'**) that has been attacked less than 10 times in the whole period from Jan 2009 to March 2013. This implies huge target availability for cybercriminals. It could be attributed to huge amounts of trial and error and R&D by the cybercriminals. Or, the attacks are executed very selectively for specific social/political/personal reasons.
- ❖ The number of active botnets had a decreasing trend from mid-2011 onwards, which might be attributed to the Zeus take down efforts being done by national governments and security firms.
- ❖ Botnets differ in their attack strategy in regards to the geographical location of the target. The data indicated that among the 10 most active botnets, some attacked a wide variety of countries while others attacked more concentrated and only to a small number of countries with similar attributes

8.1.3 TEST OF HYPOTHESES

H1: Banks with word 'bank'

The idea about the possibility that domains with word 'bank' may be attacked more than domains with 'bank' word was primarily put forward by experts in interviews. However, the result of our analysis does not support this hypothesis. Accordingly, we conclude that attackers actually think more professionally and have more knowledge about the popular domains.

H2: Banks with English webpages

Similar to the first hypothesis, the initial idea in regards to this hypothesis came from expert interviews. Our dataset supports the hypothesis only for the attacked within the EU region in all of the test years (2012 and pooled data). We argued that, the reason might be explained by the fact that the EU region contains countries with a wide variety of languages, which have similarities in terms of infrastructure, etc. So if we consider that everything rather than language is more or less the same for attacking banks in EU region, then it is logical that attackers prefer domains with English options over other domains. This is simply because in this way they can attack a wide variety of domains in different countries in EU with less effort (in term of language knowledge).

H3: Effect of broadband penetration

As it is mentioned in by experts, broadband penetration is one of the technologies that enables of malware and botnets and our dataset partially support it as well. That is, our results illustrate that as the number of broadband users increase as they become targets that are more valuable because they can be potential targets looking from cybercriminals' perspective. However, the statement holds true only in case of the global region.

H4: Effect of GDP

We expected that cybercriminals have an intuition about the statistics in regards to some country related statistics based on our conceptual model. For instance, they more or less know which countries have high GDP. However, our measures indicate that the hypothesis only holds true for globally but not for EU. The reason boils down to the sample variation; that is GDP matters, but when they are within a specific range, than other factors are more important.

H5: Effect of infection rate

We expected that cybercriminals have an intuition about the statistics in regards to the infection rate of different countries. For instance, they more or less know which countries have high infection rate. However, the result of our analysis indicates that criminals may not make their attack decisions completely based on such statistics but

rather makes use any of their chance for executing an online attack on a wide range of target domains. The results confirm our conclusion as none of them has correlation with number of attacks. A point of caution should be raised that for this hypothesis we are implicitly assuming proximity between where end-users and end-services are located.

EU versus Global

In the second hypothesis, we saw that domains with English domain pages within the EU region have been attacked more on average comparing to the global region. We explained that, the underlying reason of the above result might be explained by the fact that the EU region contains countries with a wide variety of languages, which have similarities in terms of infrastructure, etc. This means that if we suppose attackers require more or less the same infrastructures for attacking EU region banks, apart from language, then it is logical if they attack domains with English language option more than other domains because requires less effort. However, the hypothesis does not hold true globally because the globally countries are the variety of infrastructures in different countries around the world increase and thus the pattern cannot be seen anymore.

Looking at the result of hypotheses 3 and 4 we can see this fact as well: comparing the 'coefficient of variance' or CV of the two variables in EU versus global region, we realized that variety is more in both of the variables within the global region. Likewise, considering the fact that 'GDP' and 'broadband users' are infrastructure variables, it is logical that attackers take these variables into account, or these variables become more important where variety is more comparing to the EU region where variety is less.

The result of our decision-tree analysis also confirms the above statement and the result of our hypotheses in bivariate analysis. In the visualized tree being discussed in chapter 7, we saw that in EU region variable 'with English' was the most important results while in others it was not the case.

8.1.4 SUMMARY

Recalling the research question mentioned in the first section of this chapter we would conclude that:

Yes, it is possible to extract intelligence on criminal attack patterns and target selection from the configuration files which financial malware use as instructions for its operations. Taking the case study of Zeus malware, 'Average number of botnets attacking a domain (per week) is the best metric that is able to express popularity of targets among cybercriminals. Among independent variables investigated, broadband penetration and GDP showed effect on cybercriminals' decision in regards to selection of country only in the global region. It is also proved that domains with English option have been selected more often by cybercriminals within the EU region.

8.2 RECOMMENDATIONS AND PRACTICAL IMPLICATIONS

8.2.1 SCIENTIFIC RECOMMENDATIONS

- In this research, we conceptualized the problem of online banking fraud using Routine Activity Theory (RAT) from criminology literature. Although in his research, Yar (2005) studies the extent to which RAT is exposable to virtual crime, no other literature has been found that a real-world problem in the virtual space gets investigated using RAT. Therefore, this research can be used as a starting point for studying different types of virtual crime using frameworks available in the non-virtual criminology literature such as RAT. These findings can be further elaborated in other research in the field of cybersecurity in future.
- Instead of simply counting the number of times a domain is seen in the configuration files of a malware (we study Zeus financial malware in our research), we developed more advanced methods to extract meaningful results in regards to the 'number of times a domain is attacked'. In addition, some of the

variables introduced in this research, such as *'changes in defense activity level per bank'* can be used to potentially investigate cybercriminals' activities in reaction to the changes of defender's defense measures. Due to the limitations of a master's thesis, here, only one method is tested leaving room for further research based on other introduced variables.

- In this research, based on our measurement model (presented in chapter six), five empirical hypotheses are tested. The limited number of tested hypotheses can primarily be explained due to the scope of this research and secondly, according to our limited access to empirical data concerning external and explanatory factors. The dependent variable that is statistically tested in this research can further be a subject for analysis using more explanatory factors as well as other methods of analysis such as different data mining methods (we used decision-tree analysis in this research).

8.2.2 RECOMMENDATIONS FOR FOX-IT

In theory, the dataset provided by Fox-IT can be improved by:

- I. Including more metadata about the stored data (could be usable by researchers from different fields for instance). The metadata provided for this research was not complete and does not include all of the required fields. More documentation can also help researchers to be more efficient in term of understating the data.
- II. Creating uniform output for their decryption tool or updating the whole dataset using the last version of the decryption tool where their tool is changing over time. In the dataset we used in this research, the configuration files were decrypted with three-four different formats.
- III. Storing the data in a more structured way (i.e. all information in My-SQL tables and no text files)

We are aware of the fact that in practice applying all these recommendations may not be simple. However, some other actions would be helpful for that matter. The process of future research within Fox-IT can be improved if the number of experts that have information about a dataset would be more than one person. This would facilitate the research process and avoid possible bottlenecks that may occur during the research process due to the dependency on the knowledge of one expert. Moreover, for ordering, filtering and aggregating their datasets, Fox-IT can define research projects for programmers and social science students as a collaborative work. The contribution of this work (the already built dataset) can also be used as input in other future research in Fox-IT.

8.3 DISCUSSION

In the initial parts of this report, we argued that during the past decades, most of the problems in the field of cybersecurity and more specifically online banking fraud have been investigated from technological perspective. That is, these problems were considered solely as technical problems which required technical security solutions. However, as stated by (Anderson et al., 2012) most of the problems in the field of cybersecurity are more economic than technical. Based on their explanation, economics introduces breakthroughs in understanding problems of security, based on explaining the actors involved and their socio-economic incentives.

Accordingly, in this research we have reviewed Economics of Information Security literature to study the possible explanatory factors that may explain why certain targets are by cybercriminals. At the end of our literature review we realized that Economics of information security is able to provide insights to some of the economic incentives of cybercriminals in selecting their targets, as well as their economic behavior in imitating targets of others as a result of '*informational cascade*'.

However, for finding the factors that may influence the target selection by cybercriminals, we needed a field that would be able to more systematically explain the online banking crime as a dependent variable, along with the independent or explanatory factors that influence the occurrence of the crime . Consequently, we reviewed criminology literature as one of the most relevant fields to cybercrime. Using the Routine Activity Theory from criminology literature and further developing the factors identified in the economics literature, we build a conceptual framework that we expect to be able to explain situations of online banking attacks and the factors that influence attack decisions of cybercriminals.

However, after conducting the analysis in chapter 7 with the variables discussed in the empirical model in chapter 6, we realized that among the relationships that are tested, many did not hold. It should be reminded that most of the explanatory factors include in the measurement model were suggested by the field experts. Potentially two point of discussion could be raised based on the acquired results:

- ❖ RAT is being used to put the independent variables suggested by experts into work. Therefore, we can understand that even for the people in the frontline it is hard to see the patterns in this field, taking of course into account that our access to secondary empirical data about other suggested variables was limited.
- ❖ It is intriguing that the majority of the relationships that we identified based on the theory do not hold and this questions the framework we developed based in RAT. It can be recognized that the intermediate variables suggested by RAT theory, do not seem to be very important drivers here. That of course, does not mean the theory is wrong but it does suggest that one should not presume that the RAT theory works in regards to the problem of online banking fraud in cyberspace as well as it works outside the cyberspace for non-virtual crimes.

8.4 LIMITATIONS AND SUGGESTIONS FOR FURTHER RESEARCH

It should be noted that, a research is not complete unless it includes discussions on reliability, validity and limitations of the work along with the suggestions for future research. The issues above will be addressed in this final section of the thesis. The limitations that are going to be discussed in this section will mainly address the question that whether we have adequately answered our research question.

8.4.1 RELIABILITY OF INSTRUMENTS

The reliability of a measurement instrument determines the extent to which a measurement instrument is accurate (Velde, Jansen, & Anderson, 2007). One such measure for virtual space research is repeatability. Because the empirical data of this research is provided by a third party, the issue of repeatability does not hold true and thus, we will end up with the same dataset no matter of the number of times we run our scripts. Reliability of the Zeus malware data itself could be the subject of further checks but because the Zeus malware data has collected at a certain time in the past, repeatability again cannot take place.

8.4.2 VALIDITY

The notion of validity mainly points to two issues: the extent to which the research instrument measures what is aimed to be measured (Validity of Instruments) as well as the extent to which a research strategy results a type of the conclusions that we draw from it (validity of research strategy and results) (Velde et al., 2007).

Validity of Instruments

According to the first point, the validity of our measurement instruments are addressed in depth in chapter five where we introduced the limitations of each of the intermediate variables that helped explain target popularity in our dataset. It should be noted that limitations in measurement instruments are common among different types of research. Table 20 contains the summary of all the point being observed in this research.

Validity of the Research Strategy and Results

According to the second point, the validity of research results depends on several factors and validity of sample is one of those. The validity of our sample data is checked in appendix one, in data triangulation section. Although our sample did not match with other samples published in security reports in terms of top attacked domains (which can be attributed to the way different firms counts the number of attacked domains in the configuration files), it was representative in terms of number of configuration files sent by Zeus malware samples. However, it should be noted that our sample covered a time span of 4 years that was more than all of the sources with which we have crosschecked our data. Therefore, it is likely that our data provide more precise and reliable information about Zeus samples than other sources.

Moreover, in terms of validity of results, the research concerns with the issue of our limited access to independent variables that were themselves, mostly proxies for other variables. This ultimately would increase the risk of creation of errors or biases in the measurement and results.

TABLE 20-LIMITATIONS OF THE RESEARCH

Limitation	Description
Validity of measurement: extracted target domains	<ul style="list-style-type: none"> • Limitations of our automated scripts in dealing with exceptions in obfuscated URLs. • Limitations in categorizing domains (i.e. same domain names with different TLDs, different web pages of a same bank).
Validity of measurement: extracted infection time	<ul style="list-style-type: none"> • Our scripts set the first time each configuration file is seen as the infection time, since we have no information about whether the actual infection has happened or not.
Validity of measurement: different login keys, different botnets	<ul style="list-style-type: none"> • Limitation of our assumption about setting each login key as a unique botnet while a same botnet may use different login keys (due to different reason like change in the login key). • Limitation of our assumption in considering each login key as a unique botnet while different botnets may use same login keys. • Limitations in data (i.e. No login key was found for a few number of configuration files).
Validity of measurement: the dependent variable	<ul style="list-style-type: none"> • Limitation in our assumption about the number of times attackers actually perform their attacks by sending configuration files in a week (i.e. we assumed that attackers only sent one configuration file per week).
Validity of measurement: independent variables	<ul style="list-style-type: none"> • Most of our independent variables come from secondary sources, such as the World Bank or ITU, asserting their validity. So there should exist no particular issues in this regard.
Validity of results: sample validity	<ul style="list-style-type: none"> • Our Zeus malware sample partially seems to be a representative of worldwide malware traffic (see appendix one).
Validity of results: content validity	<ul style="list-style-type: none"> • Limitation in our analysis due to excluding the score of dependent and independent variables in regards to the countries that were not in our dataset. • Limitation in the number of independent variables for testing the method that is introduced in chapter 5.
Validity of results : construct validity	<ul style="list-style-type: none"> • Limitation resulted by not fully examining the interactions between independent variables. The fallacy of ‘relation does not imply causation’ might occur in certain cases. • Conceptual framework does not incorporate dynamic effects. • Conceptual framework does not fully investigate the crime in a specific time and place (as it is one of the basics of RAT).

8.4.3 RECOMMENDATIONS FOR FURTHER RESEARCH

This research can be extended in future in several ways. Moreover, there are some points that can be improved as for the future research. Since Fox-IT provides the data of the research, the suggestions in Table 21 in regards to the data part can be useful for the company's future strategy for data gathering/aggregating.

TABLE 21-SUGGESTIONS FOR FURTHER RESEARCH

Suggestion	Description
Improving the conceptual model/Future use of RAT	<ul style="list-style-type: none"> • The conceptual model can be improved in terms of dynamicity. • The customized RAT model can be further developed and used for analysing other topics in cybersecurity field.
Analysing other introduced dependent variables	<ul style="list-style-type: none"> • Analysing popularity of target from different perspectives (other dependent variables introduced in chapter 5, e.g. the number of weeks a domain is attacked). • The same dependent variable can be the subject of future studies taking into account the number of botnets that are active each week. • The impact of different authentication methods can be extracted in future studies from configuration files.
Enriching independent variables	<ul style="list-style-type: none"> • Adding more reliable data in respect to the online behaviour of users on the Internet (e.g. online shopping behaviour, online banking behaviour, browsing time etc.) • Adding more data in regards to the banking sector (e.g. banking policies, cross-national money transferring policies etc.) • Considering other external factors such as regulation as independent variable.
Data analysis	<ul style="list-style-type: none"> • Use of other data mining methods such as cluster analysis, factor analysis, regression analysis

Literature List

- AhnLab. 2012. Malware Analysis: Citadel: AhnLab ASEC (AhnLab Security Emergency response Center).
- Akerlof, G. A. 1970. The market for "lemons": Quality uncertainty and the market mechanism. *The quarterly journal of economics*: 488-500.
- Alaganandam, H., Mittal, P., Singh, A., & Fleizach, C. 2007. Cybercriminal Activity. **Accessed on 14th November**.
- Alexa. 2013. **Alexa** The Web Information Company. viewed on July 2013 from <<http://www.alexacom/topsites>>
- Amit, R., & Zott, C. 2001. Value creation in e-business. *Strategic management journal*, 22(6-7): 493-520.
- Anderson, R. 2001. **Why information security is hard-an economic perspective**. Paper presented at the Computer Security Applications Conference, 2001. ACSAC 2001. Proceedings 17th Annual.
- Anderson, R., Barton, C., Böhme, R., Clayton, R., van Eeten, M. J. G., Levi, M., Moore, T., & Savage, S. 2012. Measuring the cost of cybercrime.
- Anderson, R., Moore, T., Nagaraja, S., & Ozment, A. 2007. Incentives and information security, Vol. 881339837: ISBN-13.
- Asghari, H. 2010. **Botnet mitigation and the role of ISPs: A quantitative study into the role and incentives of Internet Service Providers in combating botnet propagation and activity**. Delft University of Technology.
- Bauer, J. M., & Van Eeten, M. J. 2009. Cybersecurity: Stakeholder incentives, externalities, and policy options. *Telecommunications Policy*, 33(10): 706-719.
- Bikhchandani, S., Hirshleifer, D., & Welch, I. 1992. A theory of fads, fashion, custom, and cultural change as informational cascades. *Journal of political Economy*: 992-1026.
- Böhme, R., & Moore, T. 2009. The Iterated Weakest Link--A Model of Adaptive Security Investment.
- Calisir, F., & Gumussoy, C. A. 2008. Internet banking versus other banking channels: Young consumers' view. *International Journal of Information Management*, 28(3): 215-221.
- Castells, M. 2003. **The Internet galaxy: Reflections on the Internet, business, and society**. New York: Oxford University Press.
- Choo, K.-K. R. 2011. The cyber threat landscape: Challenges and future research directions. *Computers & Security*, 30(8): 719-731.
- Claessens, J., Dem, V., De Cock, D., Preneel, B., & Vandewalle, J. 2002. On the security of today's online electronic banking systems. *Computers & Security*, 21(3): 253-265.
- Cohen, L. E., & Felson, M. 1979. Social change and crime rate trends: A routine activity approach. *American sociological review*: 588-608.
- De Bruijn, H., & Ten Heuvelhof, E. 2012. **Management in Networks: On multi-actor decision making**: Routledge.
- Denning, D. 2000. Information warfare and security .Vol. 4. Reading MA: Addison-Wesley.
- Douglas, T., & Loader, B. D. 2000. **Cybercrime: Security and surveillance in the information age**: Routledge.
- F-Secure. 2012a. **Threat Report H1 2012**. accessed on 2013.
- F-Secure. 2012b. **Threat Report H2 2012**. accessed on 2013.
- Florêncio, D., & Herley, C. 2010. **Phishing and money mules**. In Information Forensics and Security (WIFS), 2010 IEEE International Workshop on (pp. 1-5). IEEE.
- Florêncio, D., & Herley, C. 2011. Where Do All The Attacks Go? *Economics of Information Security and Privacy III* (pp. 13-33). Springer New York.
- Goertz, G., & Mahoney, J. 2005. Two-level theories and fuzzy-set analysis. *Sociological Methods & Research*, 33(4): 497-538.
- Hammock, M. 2010. A Review of the Economics of Information Security Literature. **A Review of the Economics of Information Security Literature (August 15, 2010)**.
- Hussain, A., Heidemann, J., & Papadopoulos, C. 2003. **A framework for classifying denial of service attacks**. Paper presented at the Proceedings of the 2003 conference on Applications, technologies, architectures, and protocols for computer communications.
- Hutchinson, D., & Warren, M. 2003. Security for internet banking: a framework. *Logistics Information Management*, 16(1): 64-73.

- Jakobsson, M., & Myers, S. 2006. *Phishing and countermeasures: understanding the increasing problem of electronic identity theft*: Wiley-Interscience.
- Jaleshgar, R. 1999. Document trading online. *Information Week*, 755: 136.
- Kothari, C. 2009. *Research methodology: methods and techniques*: New Age International.
- Krebs, B. 2011. *Zeus Innovations: 'No-\$H!+ Reports'*. Krebs on Security. viewed on 2013 from <<http://krebsonsecurity.com/2011/03/zeus-innovations-no-h-reports/>>
- Macdonald, D. 2011. *Zeus: God of DIY Botnets*. FortiGuard Center. viewed on May 2013 from <<http://www.fortiguard.com/legacy/analysis/zeusanalysis.html#3>>
- Mannan, M., & van Oorschot, P. C. 2008. *Security and usability: the gap in real-world online banking*. Paper presented at the Proceedings of the 2007 Workshop on New Security Paradigms.
- Marcus, D., & Sherstobitoff, R. 2012. *Dissecting Operation High Roller*. accessed on May 2013.
- McAfee. 2012. *McAfee Threats Report: Third Quarter 2012*. M. Labs, accessed on June 2013.
- McCullagh, A., & Caelli, W. 2005. *Who goes there? Internet banking: A matter of risk and reward*. Paper presented at the Information Security and Privacy.
- Moore, T., & Anderson, R. 2011. *Economics and Internet security: A survey of recent analytical, empirical and behavioral research*. Harvard University Computer Science Group.
- Moore, T., & Clayton, R. 2007. *Examining the impact of website take-down on phishing*. Paper presented at the Proceedings of the anti-phishing working groups 2nd annual eCrime researchers summit.
- Moore, T., Clayton, R., & Anderson, R. 2009. The economics of online crime. *The Journal of Economic Perspectives*, 23(3): 3-20.
- Nevens, T. M. 1999. The mouse that roared. *McKinsey Quarterly*: 145-148.
- NYTimes. 2013. *A National Priority and a Business Priority*. viewed on March 2013 from <<http://www.nytimes.com/roomfordebate/2013/02/21/should-companies-tell-us-when-they-get-hacked/countering-security-breaches-is-a-national-priority-and-a-business-priority>>
- OECD. 2007. *Malicious Software (Malware): A Security Threat to the Internet Economy*. accessed on June 2013.
- Olson, M. 2009. *The logic of collective action: public goods and the theory of groups*: Harvard University Press.
- Premchaiswadi, N., Williams, J. G., & Premchaiswadi, W. 2009. A Study of an On-Line Credit Card Payment Processing and Fraud Prevention for e-Business. In T. Bastiaens, J. Dron, & C. Xin (Eds.), *World Conference on E-Learning in Corporate, Government, Healthcare, and Higher Education 2009*: 2199-2206. Vancouver, Canada: AACE.
- Raafat, R. M., Chater, N., & Frith, C. 2009. Herding in humans. *Trends in cognitive sciences*, 13(10): 420-428.
- Rittel, H. W., & Webber, M. M. 1973. Dilemmas in a general theory of planning. *Policy sciences*, 4(2): 155-169.
- RSA. 2012. CITADEL TROJAN OUTGROWING ITS ZEUS ORIGINS, *Fraud Report*: RSA.
- RSALabs. 2012. *New Trojan Ice IX Written Over Zeus' Ruins*. RSA. viewed on from <<https://blogs.rsa.com/new-trojan-ice-ix-written-over-zeus-ruins/>>
- Schneier, B. 2011. *Secrets and lies: digital security in a networked world*: Wiley.
- Sherstobitoff, R. 2013. Inside the World of the Citadel Trojan McAfee Labs. accessed on June 2013
- Shields, P. M., & Tajalli, H. 2006. Intermediate theory: The missing link in successful student scholarship. *Journal of Public Affairs Education*: 313-334.
- Simon, H. A. 1991. Bounded rationality and organizational learning. *Organization science*, 2(1): 125-134.
- Stone-Gross, B. 2012. *The Lifecycle of Peer-to-Peer (Gameover) Zeus*. Dell SecureWorks Counter Threat Unit(TM) Threat Intelligence. viewed on 2013 from <http://www.secureworks.com/cyber-threat-intelligence/threats/The_Lifecycle_of_Peer_to_Peer_Gameover_Zeus/>
- Symantec. 2012. *Zeusbot/Spyeye P2P Updated, Fortifying the Botnet*. Symantec. viewed on 2013 from <<http://www.symantec.com/connect/blogs/zeusbotspyeye-p2p-updated-fortifying-botnet>>
- Symantec. 2013. *Internet Security Treath Report 2013*. accessed on June 2013.
- TrendMicro. 2010. *Zeus and Its Continuing Drive Towards Stealing Online Data*. accessed on April 2013.
- TrendMicro. 2013. *Security Threats to Business, the Digital Lifestyle, and the Cloud*. accessed on June 2013.
- Utakrit, N. 2009. Review of Browser Extensions, a Man-in-the-Browser Phishing Techniques Targeting Bank Customers.
- Van Eeten, M., Bauer, J., Asghari, H., Tabatabaie, S., & Rand, D. 2010. *The role of internet service providers in botnet mitigation an empirical analysis based on spam data*. TPRC

- Van Eeten, M. J. G., & Bauer, J. M. 2008. *Economics of malware: Security decisions, incentives and externalities*. No. 2008/1. OECD Publishing.
- Velde, v. d., Jansen, P. G., & Anderson, N. 2007. *Guide to management research methods*: Blackwell.
- Vrancianu, M., & Popa, L. A. 2010. Considerations Regarding the Security and Protection of E-Banking Services Consumers' Interests. *The Amfiteatru Economic Journal*, 12(28): 388-403.
- Wall, D. 2001. 1 Cybercrimes and the Internet. *Crime and the Internet*: 1.
- Witten, I. H., & Frank, E. 2005. *Data Mining: Practical machine learning tools and techniques*: Morgan Kaufmann.
- Wyke, J. 2011. *What is Zeus?* SOPOHS, accessed on June 2013.
- Yar, M. 2005. The Novelty of 'Cybercrime' An Assessment in Light of Routine Activity Theory. *European Journal of Criminology*, 2(4): 407-427.

Appendix One - Data Triangulation

One of the most important steps in an empirical research is to realize whether the data being used in the research is an 'indicative sample' of the actual population. Accordingly, in this research we have to check to see whether the Fox-IT Zeus malware dataset could be an indicative sample of online banking malware attacks. For this, we have to compare our data with the already available statistics in regard to online banking malware attacks being published by security firms or banking industry.

TABLE 22-LIST OF THE WELL-KNOWN AVAILABLE SECURITY REPORTS AND BLOGS

Security Firm	Report Name
AhnLab	Malware Analysis: Citadel
F-Secure	Threat Report H1
McAfee	Threats Report
Microsoft	Security Intelligence report
RSA	Fraud Report
Sophos	Security Threat Report
Symantec	Intelligence Report
TrendMicro	Trend Micro annual threat roundup and forecast
Abuse.ch	Zeus Tracker
Alexa.com	The web information company

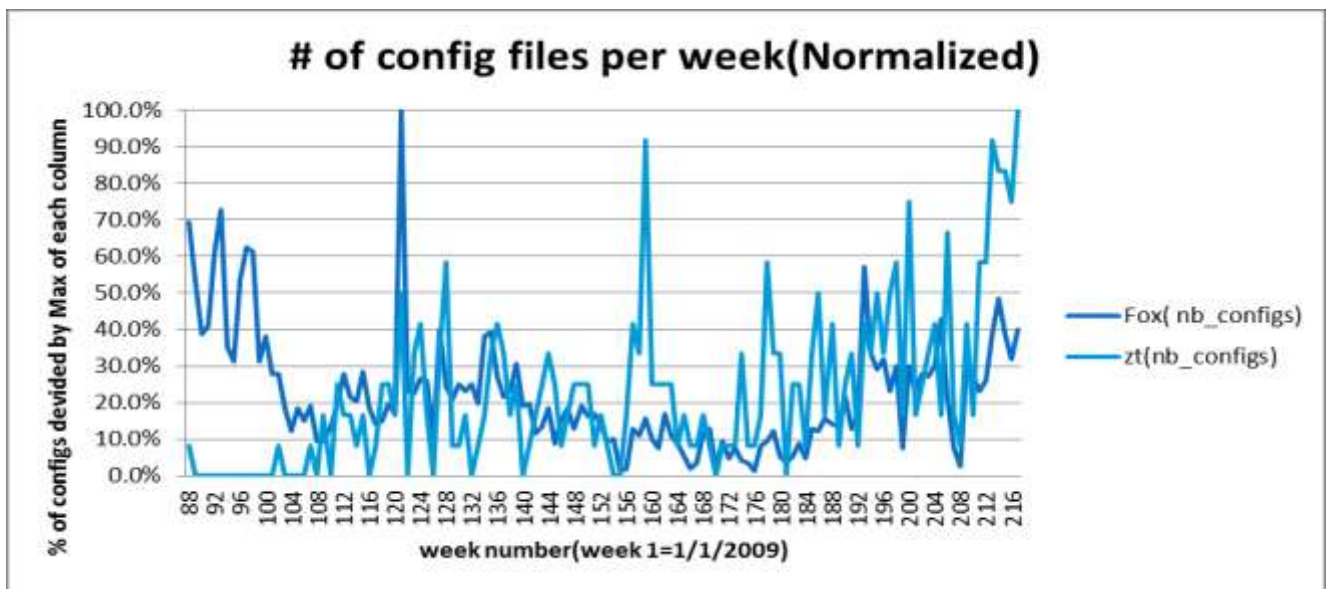
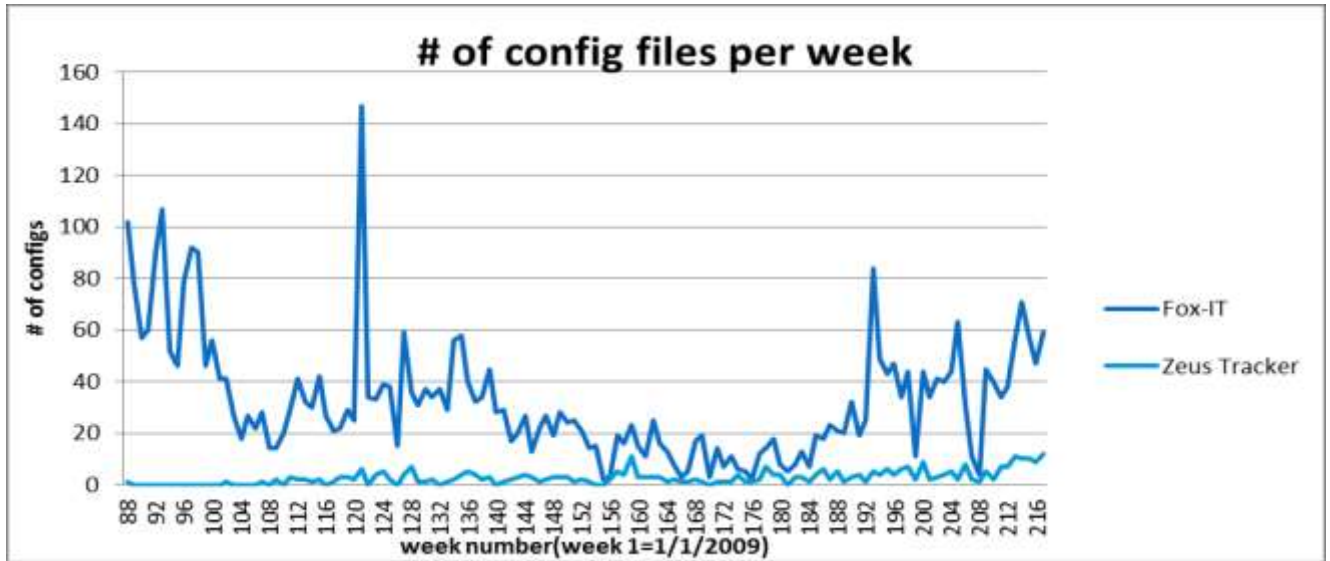
The choice of information to be used from various existing industry and security reports is highly dependent to the relevancy of points that are reported in the industry report with our dataset. Among all of the available choices we have to select the ones that could be a correct point of comparison with our dataset. As it mentioned previously, due to confidentiality reasons, most of the security reports do not publish a lot of statistics in regards to specific targets of online financial malware attacks and their specifications. This leaves us with relatively few choices. After reviewing the security reports mentioned in Table 22, we found the following the most relevant ones to the information existed in our dataset:

- The number of Zeus configuration files /host received by Zeus Tracker
- The list of targeted domains Zeus and SpyEye Trojan
- The list of targeted countries by Zeus and its variants
- The list of most visited domains (Alexa Rankings)

In the following section we will explain the above two points more in details.

COMPARISON BETWEEN NUMBER OF CONFIGURATION FILES

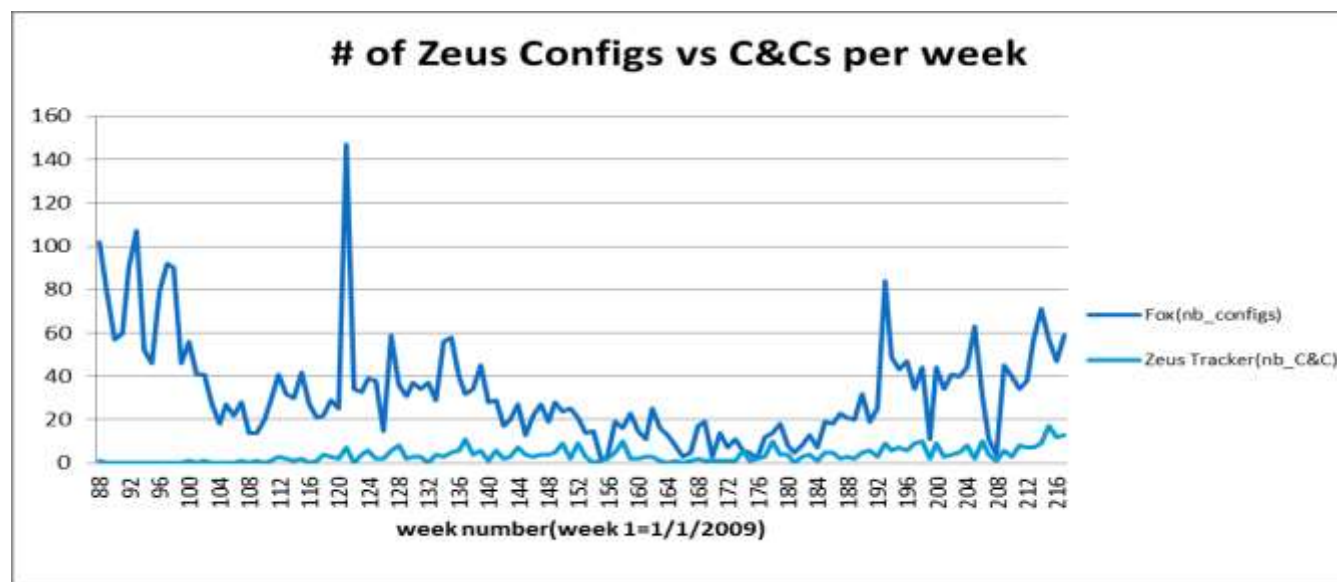
The first triangulation attempt uses Abuse.ch data on number of Zeus configuration files from 2010 till 2013. The figures below display the comparison between Fox-IT and Zeus Tracker datasets. The first two figures contain difference between numbers of configuration files per week in the two datasets. The Zeus tracker graph is built upon the list and date of configuration files that has been tracked by the Zeus tracker. The second graph is the normalized version of the first graph, being normalized through dividing the records of each group by the maximum of the group.



As it can be seen from the above figure the trend of the two datasets is more different in mid-2010 (weeks 88-100) while in other times they more or less follow a same trend. Also as it is more obvious from the first figure, the number of configuration files in the dataset provided by Fox-IT is much more than Zeus tracker dataset. This can increase the probability that our Fox-IT sample be more representative than Zeus tracker's.

COMPARISON BETWEEN NUMBER OF CONFIGURATION FILES VERSUS NUMBER OF COMMAND AND CONTROL (C&C) FILES

The second attempt of triangulation is being done by comparing the number of configuration files extracted from Fox-IT dataset with the number of C&C servers tracked by Zeus tracker listed in Abuse.ch website. Same as the previous graphs, the second graphs in this section is also the normalized version of the first graph. Although the points of comparison are not exactly identical, but we can argue that number of Zeus configuration files in a period could be a factor of number of command and control servers in the same period.



As the second figure illustrates, the number of Fox-IT's Zeus configuration files is a more or less identical proportion of Zeus tracker's number of Zeus C&C files except from the first 20 weeks (88-108).

COMPARISON BETWEEN TOP-20 BANKS TARGETED BY ZEUS VERSUS SPYEYE TROJAN

The figures below display the difference between top-20 attacked banks by Zeus from Fox-IT data and SpyEye from F-secure dataset in the period of Jun2011 to may2012 (F-Secure, 2012a). It should be mentioned that, the ranks are only based on raw counts of the number domain names seen in the configuration files conducted only for the purpose of data triangulation. As it is obvious from the results, although some of targeted domains of both Trojans are the same, in general top domains and their associated orders are different.

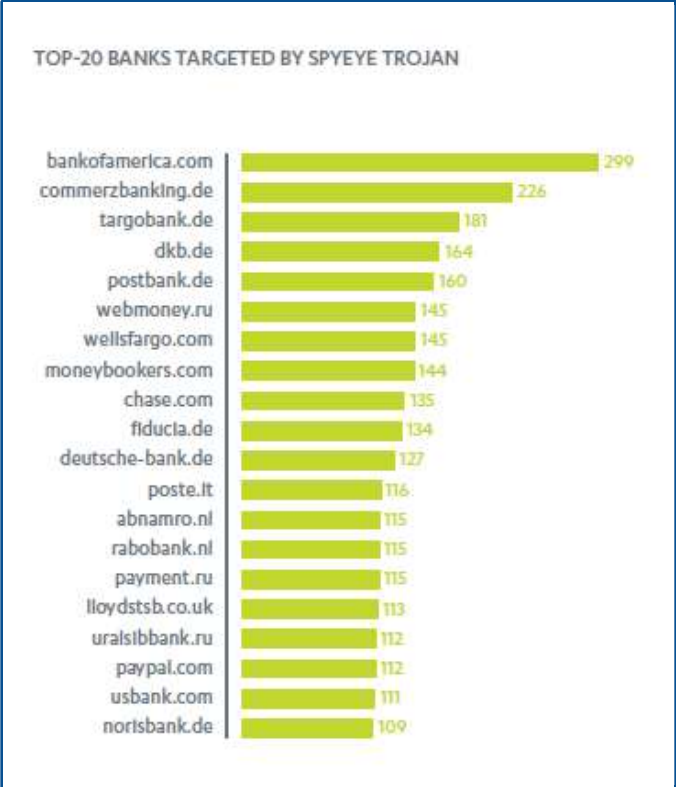
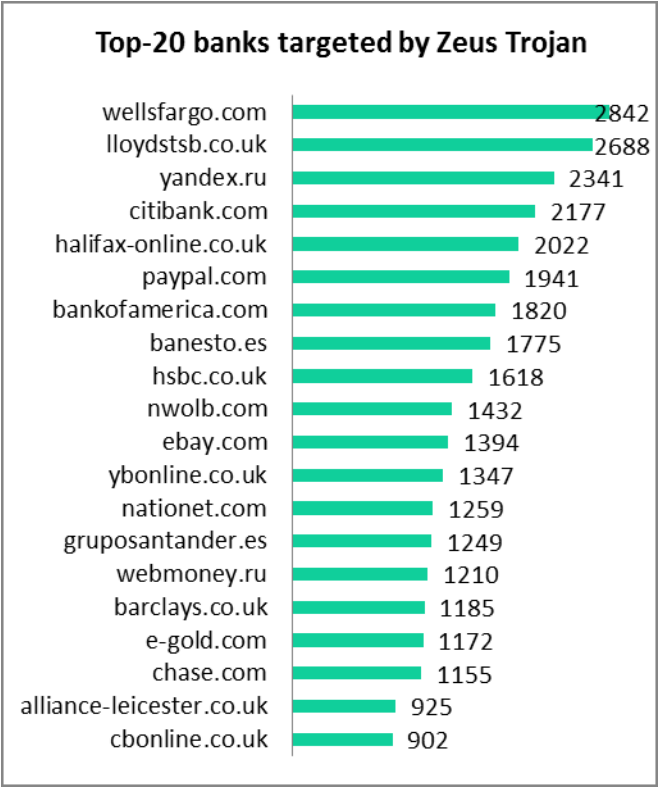


FIGURE 86-FOX-IT (LEFT) VERSUS F-SECURE (RIGHT)

COMPARISON BETWEEN TOP ATTACKED COUNTRIES

The attacked countries can also be a point for cross-checking Fox-IT dataset against datasets from other sources. F-secure threat report H2 2012 contains data about top countries attacked by P2P Zeus Trojan from late August till mid November 2012. It should be noted that because in the Fox-IT dataset no distinction is made between variants of Zeus, the below graphs include the top domains attacked by all variants of Zeus in that period not only the P2P version. Accordingly, comparing the two graphs is a bit hard. However, as it is obvious from the below figure, the top attacked countries in the two data sets does not completely follow the same order which can partially be attributed to the holistivity of Fox-IT dataset comparing to F-Secure P2P dataset.

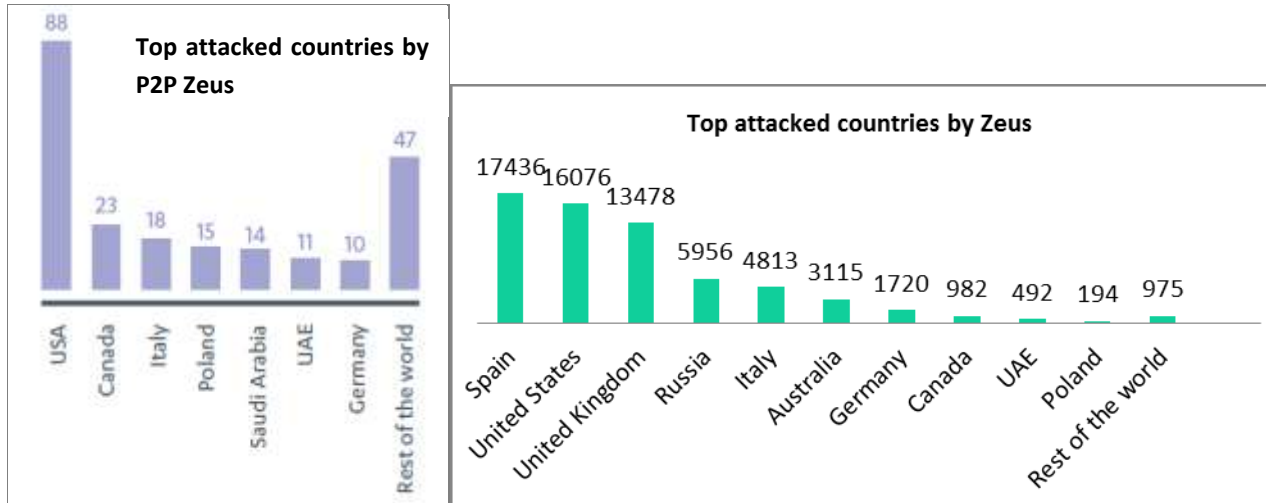


FIGURE 87-F-SECURE P2P ZEUS TOP ATTACKED COUNTRIES (LEFT) AGAINST FOX-IT ZEUS TOP ATTACKED COUNTRIES (RIGHT)(F-SECURE, 2012B)

COMPARISON BETWEEN PROPORTION OF EMAILS CONTAIN URL MALWARE AND NUMBER OF ZEUS CONFIGURATION FILES

Another existing online threat group that is to some extent comparable with our data set is the statistics in regards to the number of emails contain URL Malware. Although number of configuration files sent everyday by C&C servers should not be identical to the number of email contain URL malware, we expect a kind of similar trend in the both groups, certainly because being infected by a malware is as of initial steps of an online banking attack.

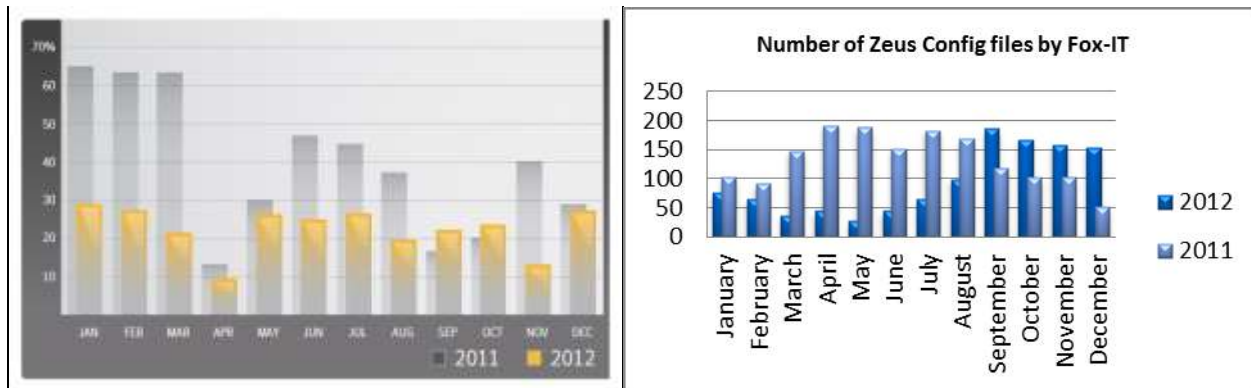


FIGURE 88-PROPORTION OF EMAIL TRAFFIC CONTAINING URL MALWARE FROM SYMANTEC DATA (LEFT) AGAINST FOX-IT NUMBER OF ZEUS CONFIGURATION FILES(RIGHT)(SYMANTEC, 2013)

As it can be determined from the above figures, although the absolute numbers are different, the overall trend of infections are identical expect for the last months of 2012. For instance, in Fox-It data number of configuration files in November and December 2012 is higher than 2011 while it is not the case in the data reported by Symantec.