

Data Mining and Machine Learning in Cybersecurity

Sumeet Dua and Xian Du



 CRC Press
Taylor & Francis Group
AN AUERBACH BOOK

Data Mining and Machine Learning in Cybersecurity

Sumeet Dua and Xian Du



CRC Press

Taylor & Francis Group

Boca Raton London New York

CRC Press is an imprint of the
Taylor & Francis Group, an **informa** business
AN AUERBACH BOOK

Auerbach Publications
Taylor & Francis Group
6000 Broken Sound Parkway NW, Suite 300
Boca Raton, FL 33487-2742

© 2011 by Taylor and Francis Group, LLC
Auerbach Publications is an imprint of Taylor & Francis Group, an Informa business

No claim to original U.S. Government works

Printed in the United States of America on acid-free paper
10 9 8 7 6 5 4 3 2 1

International Standard Book Number-13: 978-1-4398-3943-0 (Ebook-PDF)

This book contains information obtained from authentic and highly regarded sources. Reasonable efforts have been made to publish reliable data and information, but the author and publisher cannot assume responsibility for the validity of all materials or the consequences of their use. The authors and publishers have attempted to trace the copyright holders of all material reproduced in this publication and apologize to copyright holders if permission to publish in this form has not been obtained. If any copyright material has not been acknowledged please write and let us know so we may rectify in any future reprint.

Except as permitted under U.S. Copyright Law, no part of this book may be reprinted, reproduced, transmitted, or utilized in any form by any electronic, mechanical, or other means, now known or hereafter invented, including photocopying, microfilming, and recording, or in any information storage or retrieval system, without written permission from the publishers.

For permission to photocopy or use material electronically from this work, please access www.copyright.com (<http://www.copyright.com/>) or contact the Copyright Clearance Center, Inc. (CCC), 222 Rosewood Drive, Danvers, MA 01923, 978-750-8400. CCC is a not-for-profit organization that provides licenses and registration for a variety of users. For organizations that have been granted a photocopy license by the CCC, a separate system of payment has been arranged.

Trademark Notice: Product or corporate names may be trademarks or registered trademarks, and are used only for identification and explanation without intent to infringe.

Visit the Taylor & Francis Web site at
<http://www.taylorandfrancis.com>

and the Auerbach Web site at
<http://www.auerbach-publications.com>

Contents

List of Figures	xi
List of Tables	xv
Preface.....	xvii
Authors.....	xxi
1 Introduction.....	1
1.1 Cybersecurity	2
1.2 Data Mining.....	5
1.3 Machine Learning.....	7
1.4 Review of Cybersecurity Solutions.....	8
1.4.1 Proactive Security Solutions.....	8
1.4.2 Reactive Security Solutions.....	9
1.4.2.1 Misuse/Signature Detection	10
1.4.2.2 Anomaly Detection	10
1.4.2.3 Hybrid Detection	13
1.4.2.4 Scan Detection	13
1.4.2.5 Profiling Modules.....	13
1.5 Summary.....	14
1.6 Further Reading	15
References.....	16
2 Classical Machine-Learning Paradigms for Data Mining	23
2.1 Machine Learning	24
2.1.1 Fundamentals of Supervised Machine-Learning Methods	24
2.1.1.1 Association Rule Classification	24
2.1.1.2 Artificial Neural Network	25

- 2.1.1.3 Support Vector Machines27
- 2.1.1.4 Decision Trees29
- 2.1.1.5 Bayesian Network.....30
- 2.1.1.6 Hidden Markov Model.....31
- 2.1.1.7 Kalman Filter 34
- 2.1.1.8 Bootstrap, Bagging, and AdaBoost..... 34
- 2.1.1.9 Random Forest37
- 2.1.2 Popular Unsupervised Machine-Learning Methods38
 - 2.1.2.1 *k*-Means Clustering38
 - 2.1.2.2 Expectation Maximum.....38
 - 2.1.2.3 *k*-Nearest Neighbor 40
 - 2.1.2.4 SOM ANN41
 - 2.1.2.5 Principal Components Analysis41
 - 2.1.2.6 Subspace Clustering.....43
- 2.2 Improvements on Machine-Learning Methods..... 44
 - 2.2.1 New Machine-Learning Algorithms..... 44
 - 2.2.2 Resampling..... 46
 - 2.2.3 Feature Selection Methods 46
 - 2.2.4 Evaluation Methods.....47
 - 2.2.5 Cross Validation49
- 2.3 Challenges50
 - 2.3.1 Challenges in Data Mining50
 - 2.3.1.1 Modeling Large-Scale Networks50
 - 2.3.1.2 Discovery of Threats.....50
 - 2.3.1.3 Network Dynamics and Cyber Attacks51
 - 2.3.1.4 Privacy Preservation in Data Mining.....51
 - 2.3.2 Challenges in Machine Learning (Supervised Learning and Unsupervised Learning)51
 - 2.3.2.1 Online Learning Methods for Dynamic Modeling of Network Data52
 - 2.3.2.2 Modeling Data with Skewed Class Distributions to Handle Rare Event Detection52
 - 2.3.2.3 Feature Extraction for Data with Evolving Characteristics53
- 2.4 Research Directions.....53
 - 2.4.1 Understanding the Fundamental Problems of Machine-Learning Methods in Cybersecurity54
 - 2.4.2 Incremental Learning in Cyberinfrastructures.....54
 - 2.4.3 Feature Selection/Extraction for Data with Evolving Characteristics54
 - 2.4.4 Privacy-Preserving Data Mining.....55
- 2.5 Summary.....55
- References.....55

3	Supervised Learning for Misuse/Signature Detection	57
3.1	Misuse/Signature Detection	58
3.2	Machine Learning in Misuse/Signature Detection	60
3.3	Machine-Learning Applications in Misuse Detection.....	61
3.3.1	Rule-Based Signature Analysis.....	61
3.3.1.1	Classification Using Association Rules.....	62
3.3.1.2	Fuzzy-Rule-Based	65
3.3.2	Artificial Neural Network	68
3.3.3	Support Vector Machine.....	69
3.3.4	Genetic Programming	70
3.3.5	Decision Tree and CART	73
3.3.5.1	Decision-Tree Techniques.....	74
3.3.5.2	Application of a Decision Tree in Misuse Detection	75
3.3.5.3	CART	77
3.3.6	Bayesian Network.....	79
3.3.6.1	Bayesian Network Classifier	79
3.3.6.2	Naïve Bayes	82
3.4	Summary.....	82
	References.....	82
4	Machine Learning for Anomaly Detection	85
4.1	Introduction	85
4.2	Anomaly Detection	86
4.3	Machine Learning in Anomaly Detection Systems.....	87
4.4	Machine-Learning Applications in Anomaly Detection	88
4.4.1	Rule-Based Anomaly Detection (Table 1.3, C.6).....	89
4.4.1.1	Fuzzy Rule-Based (Table 1.3, C.6)	90
4.4.2	ANN (Table 1.3, C.9).....	93
4.4.3	Support Vector Machines (Table 1.3, C.12)	94
4.4.4	Nearest Neighbor-Based Learning (Table 1.3, C.11).....	95
4.4.5	Hidden Markov Model.....	98
4.4.6	Kalman Filter	99
4.4.7	Unsupervised Anomaly Detection	100
4.4.7.1	Clustering-Based Anomaly Detection.....	101
4.4.7.2	Random Forests.....	103
4.4.7.3	Principal Component Analysis/Subspace	104
4.4.7.4	One-Class Supervised Vector Machine.....	106
4.4.8	Information Theoretic (Table 1.3, C.5).....	110
4.4.9	Other Machine-Learning Methods Applied in Anomaly Detection (Table 1.3, C.2)	110
4.5	Summary.....	111
	References.....	112

5	Machine Learning for Hybrid Detection	115
5.1	Hybrid Detection	116
5.2	Machine Learning in Hybrid Intrusion Detection Systems	118
5.3	Machine-Learning Applications in Hybrid Intrusion Detection....	119
5.3.1	Anomaly–Misuse Sequence Detection System.....	119
5.3.2	Association Rules in Audit Data Analysis and Mining (Table 1.4, D.4).....	120
5.3.3	Misuse–Anomaly Sequence Detection System.....	122
5.3.4	Parallel Detection System	128
5.3.5	Complex Mixture Detection System.....	132
5.3.6	Other Hybrid Intrusion Systems.....	134
5.4	Summary.....	135
	References.....	136
6	Machine Learning for Scan Detection	139
6.1	Scan and Scan Detection.....	140
6.2	Machine Learning in Scan Detection	142
6.3	Machine-Learning Applications in Scan Detection	143
6.4	Other Scan Techniques with Machine-Learning Methods	156
6.5	Summary.....	156
	References.....	157
7	Machine Learning for Profiling Network Traffic	159
7.1	Introduction	159
7.2	Network Traffic Profiling and Related Network Traffic Knowledge.....	160
7.3	Machine Learning and Network Traffic Profiling.....	161
7.4	Data-Mining and Machine-Learning Applications in Network Profiling	162
7.4.1	Other Profiling Methods and Applications.....	173
7.5	Summary.....	174
	References.....	175
8	Privacy-Preserving Data Mining.....	177
8.1	Privacy Preservation Techniques in PPDM.....	180
8.1.1	Notations.....	180
8.1.2	Privacy Preservation in Data Mining.....	180
8.2	Workflow of PPDM.....	184
8.2.1	Introduction of the PPDM Workflow.....	184
8.2.2	PPDM Algorithms.....	185
8.2.3	Performance Evaluation of PPDM Algorithms.....	185

8.3	Data-Mining and Machine-Learning Applications in PPDM.....	189
8.3.1	Privacy Preservation Association Rules (Table 1.1, A.4)	189
8.3.2	Privacy Preservation Decision Tree (Table 1.1, A.6).....	193
8.3.3	Privacy Preservation Bayesian Network (Table 1.1, A.2)	194
8.3.4	Privacy Preservation KNN (Table 1.1, A.7)	197
8.3.5	Privacy Preservation <i>k</i> -Means Clustering (Table 1.1, A.3).....	199
8.3.6	Other PPDM Methods.....	201
8.4	Summary.....	202
	References.....	204
9	Emerging Challenges in Cybersecurity	207
9.1	Emerging Cyber Threats.....	208
9.1.1	Threats from Malware	208
9.1.2	Threats from Botnets	209
9.1.3	Threats from Cyber Warfare.....	211
9.1.4	Threats from Mobile Communication	211
9.1.5	Cyber Crimes	212
9.2	Network Monitoring, Profiling, and Privacy Preservation.....	213
9.2.1	Privacy Preservation of Original Data.....	213
9.2.2	Privacy Preservation in the Network Traffic Monitoring and Profiling Algorithms.....	214
9.2.3	Privacy Preservation of Monitoring and Profiling Data	215
9.2.4	Regulation, Laws, and Privacy Preservation.....	215
9.2.5	Privacy Preservation, Network Monitoring, and Profiling Example: PRISM	216
9.3	Emerging Challenges in Intrusion Detection	218
9.3.1	Unifying the Current Anomaly Detection Systems	219
9.3.2	Network Traffic Anomaly Detection	219
9.3.3	Imbalanced Learning Problem and Advanced Evaluation Metrics for IDS.....	220
9.3.4	Reliable Evaluation Data Sets or Data Generation Tools.....	221
9.3.5	Privacy Issues in Network Anomaly Detection.....	222
9.4	Summary.....	222
	References.....	223

List of Figures

Figure 1.1	Conventional cybersecurity system	3
Figure 1.2	Adaptive defense system for cybersecurity	4
Figure 2.1	Example of a two-layer ANN framework.....	26
Figure 2.2	SVM classification. (a) Hyperplane in SVM. (b) Support vector in SVM.....	28
Figure 2.3	Sample structure of a decision tree	29
Figure 2.4	Bayes network with sample factored joint distribution	30
Figure 2.5	Architecture of HMM.....	31
Figure 2.6	Workflow of Kalman filter.....	35
Figure 2.7	Workflow of AdaBoost	37
Figure 2.8	KNN classification ($k = 5$).....	40
Figure 2.9	Example of PCA application in a two-dimensional Gaussian mixture data set.....	43
Figure 2.10	Confusion matrix for machine-learning performance evaluation	45
Figure 2.11	ROC curve representation	49
Figure 3.1	Misuse detection using “if–then” rules	59
Figure 3.2	Workflow of misuse/signature detection system.....	60
Figure 3.3	Workflow of a GP technique	71
Figure 3.4	Example of a decision tree	77
Figure 3.5	Example of BN and CPT	80
Figure 4.1	Workflow of anomaly detection system	88

Figure 4.2 Workflow of SVM and ANN testing.....95

Figure 4.3 Example of challenges faced by distance-based KNN methods 96

Figure 4.4 Example of neighborhood measures in density-based KNN methods97

Figure 4.5 Workflow of unsupervised anomaly detection 101

Figure 4.6 Analysis of distance inequalities in KNN and clustering.....108

Figure 5.1 Three types of hybrid detection systems. (a) Anomaly–misuse sequence detection system. (b) Misuse–anomaly sequence detection system. (c) Parallel detection system 117

Figure 5.2 The workflow of anomaly–misuse sequence detection system 119

Figure 5.3 Framework of training phase in ADAM.....121

Figure 5.4 Framework of testing phase in ADAM121

Figure 5.5 A representation of the workflow of misuse–anomaly sequence detection system that was developed by Zhang et al. (2008) 123

Figure 5.6 The workflow of misuse–anomaly detection system in Zhang et al. (2008)124

Figure 5.7 The workflow of the hybrid system designed in Hwang et al. (2007)125

Figure 5.8 The workflow in the signature generation module designed in Hwang et al. (2007)127

Figure 5.9 Workflow of parallel detection system128

Figure 5.10 Workflow of real-time NIDES.....130

Figure 5.11 (a) Misuse detection result, (b) example of histogram plot for user1 test data results, and (c) the overlapping by combining and merging the testing results of both misuse and anomaly detection systems 131

Figure 5.12 Workflow of hybrid detection system using the AdaBoost algorithm.....132

Figure 6.1 Workflow of scan detection143

Figure 6.2 Workflow of SPADE 145

Figure 6.3	Architecture of a GrIDS system for a department.....	146
Figure 6.4	Workflow of graph building and combination via rule sets.....	147
Figure 6.5	Workflow of scan detection using data mining in Simon et al. (2006)	150
Figure 6.6	Workflow of scan characterization in Muelder et al. (2007)	153
Figure 6.7	Structure of BAM.....	154
Figure 6.8	Structure of ScanVis	155
Figure 6.9	Paired comparison of scan patterns.....	155
Figure 7.1	Workflow of network traffic profiling.....	161
Figure 7.2	Workflow of NETMINE.....	163
Figure 7.3	Examples of hierarchical taxonomy in generalizing association rules. (a) Taxonomy for address. (b) Taxonomy for ports	164
Figure 7.4	Workflow of AutoFocus	166
Figure 7.5	Workflow of network traffic profiling as proposed in Xu et al. (2008)	167
Figure 7.6	Procedures of dominant state analysis.....	169
Figure 7.7	Profiling procedure in MINDS.....	171
Figure 7.8	Example of the concepts in DBSCAN	172
Figure 8.1	Example of identifying identities by connecting two data sets	178
Figure 8.2	Two data partitioning ways in PPDM: (a) horizontal and (b) vertical private data for DM	182
Figure 8.3	Workflow of SMC	183
Figure 8.4	Perturbation and reconstruction in PPDM.....	183
Figure 8.5	Workflow of PPDM	184
Figure 8.6	Workflow of privacy preservation association rules mining method.....	191
Figure 8.7	LDS and privacy breach level for the soccer data set.....	192
Figure 8.8	Partitioned data sets by feature subsets	193
Figure 8.9	Framework of privacy preservation KNN.....	197

Figure 8.10	Workflow of privacy preservation k -means in Vaidya and Clifton (2004)	199
Figure 8.11	Step 1 in permutation procedure for finding the closest cluster.....	200
Figure 8.12	Step 2 in permutation procedure for finding the closest cluster.....	200
Figure 9.1	Framework of PRISM	216

List of Tables

Table 1.1	Examples of PPDM	9
Table 1.2	Examples of Data Mining and Machine Learning for Misuse/ Signature Detection	11
Table 1.3	Examples of Data Mining and Machine Learning for Anomaly Detection	12
Table 1.4	Examples of Data Mining for Hybrid Intrusion Detection	13
Table 1.5	Examples of Data Mining for Scan Detection.....	14
Table 1.6	Examples of Data Mining for Profiling	14
Table 3.1	Example of Shell Command Data	63
Table 3.2	Examples of Association Rules for Shell Command Data	64
Table 3.3	Example of “Traffic” Connection Records	64
Table 3.4	Example of Rules and Features of Network Packets	76
Table 4.1	Users’ Normal Behaviors in Fifth Week	90
Table 4.2	Normal Similarity Scores and Anomaly Scores.....	91
Table 4.3	Data Sets Used in Lakhina et al. (2004a).....	106
Table 4.4	Parameter Settings for Clustering-Based Methods	109
Table 4.5	Parameter Settings for KNN.....	109
Table 4.6	Parameter Settings for SVM	109
Table 5.1	The Number of Training and Testing Data Types.....	134
Table 6.1	Testing Data Set Information.....	149

Table 8.1	Data Set Structure in This Chapter	180
Table 8.2	Analysis of Privacy Breaching Using Three Randomization Methods	187
Table 9.1	Top 10 Most Active Botnets in the United States in 2009.....	210

Preface

In the emerging era of Web 3.0, securing cyberspace has gradually evolved into a critical organizational and national research agenda inviting interest from a multidisciplinary scientific workforce. There are many avenues into this area, and, in recent research, machine-learning and data-mining techniques have been applied to design, develop, and improve algorithms and frameworks for cybersecurity system design. Intellectual products in this domain have appeared under various topics, including machine learning, data mining, cybersecurity, data management and modeling, and privacy preservation. Several conferences, workshops, and journals focus on the fragmented research topics in this area. However, transcendent and interdisciplinary assessment of past and current works in the field and possible paths for future research in the area are essential for consistent research and development.

This interdisciplinary assessment is especially useful for students, who typically learn cybersecurity, machine learning, and data mining in independent courses. Machine learning and data mining play significant roles in cybersecurity, especially as more challenges appear with the rapid development of information discovery techniques, such as those originating from the sheer dimensionality and heterogeneous nature of the network data, the dynamic change of threats, and the severe imbalanced classes of normal and anomalous behaviors. In this book, we attempt to combine all the above knowledge for a single advanced course.

This book surveys cybersecurity problems and state-of-the-art machine-learning and data-mining solutions that address the overarching research problems, and it is designed for students and researchers studying or working on machine learning and data mining in cybersecurity applications. The inclusion of cybersecurity in machine-learning research is important for academic research. Such an inclusion inspires fundamental research in machine learning and data mining, such as research in the subfields of imbalanced learning, feature extraction for data with evolving characteristics, and privacy-preserving data mining.

Organization

In Chapter 1, we introduce the vulnerabilities of cyberinfrastructure and the conventional approaches to cyber defense. Then, we present the vulnerabilities of these conventional cyber protection methods and introduce higher-level methodologies that use advanced machine learning and data mining to build more reliable cyber defense systems. We review the cybersecurity solutions that use machine-learning and data-mining techniques, including privacy-preservation data mining, misuse detection, anomaly detection, hybrid detection, scan detection, and profiling detection. In addition, we list a number of references that address cybersecurity issues using machine-learning and data-mining technology to help readers access the related material easily.

In Chapter 2, we introduce machine-learning paradigms and cybersecurity along with a brief overview of machine-learning formulations and the application of machine-learning methods and data mining/management in cybersecurity. We discuss challenging problems and future research directions that are possible when machine-learning methods are applied to the huge amount of temporal and unbalanced network data.

In Chapter 3, we address misuse/signature detection. We introduce fundamental knowledge, key issues, and challenges in misuse/signature detection systems, such as building efficient rule-based algorithms, feature selection for rule matching and accuracy improvement, and supervised machine-learning classification of attack patterns. We investigate several supervised learning methods in misuse detection. We explore the limitations and difficulties of using these machine-learning methods in misuse detection systems and outline possible problems, such as the inadequate ability to detect a novel attack, irregular performance for different attack types, and requirements of the intelligent feature selection. We guide readers to questions and resources that will help them learn more about the use of advanced machine-learning techniques to solve these problems.

In Chapter 4, we provide an overview of anomaly detection techniques. We investigate and classify a large number of machine-learning methods in anomaly detection. In this chapter, we briefly describe the applications of machine-learning methods in anomaly detection. We focus on the limitations and difficulties that encumber machine-learning methods in anomaly detection systems. Such problems include an inadequate ability to maintain a high detection rate and a low false-alarm rate. As anomaly detection is the most concentrative application area of machine-learning methods, we perform in-depth studies to explain the appropriate learning procedures, e.g., feature selection, in detail.

In Chapter 5, we address hybrid intrusion detection techniques. We describe how hybrid detection methods are designed and employed to detect unknown intrusions and anomaly detection with a lower false-positive rate. We categorize the hybrid intrusion detection techniques into three groups based on combinational methods. We demonstrate several machine-learning hybrids that raise detection accuracies in

the intrusion detection system, including correlation techniques, artificial neural networks, association rules, and random forest classifiers.

In Chapter 6, we address scan detection techniques using machine-learning methods. We explain the dynamics of scan attacks and focus on solving scan detection problems in applications. We provide several examples of machine-learning methods used for scan detection, including the rule-based methods, threshold random walk, association memory learning techniques, and expert knowledge-rule-based learning model. This chapter addresses the issues pertaining to the high percentage of false alarms and the evaluation of efficiency and effectiveness of scan detection.

In Chapter 7, we address machine-learning techniques for profiling network traffic. We illustrate a number of profiling modules that profile normal or anomalous behaviors in cyberinfrastructure for intrusion detection. We introduce and investigate a number of new concepts for clustering methods in intrusion detection systems, including association rules, shared nearest neighbor clustering, EM-based clustering, subspace, and informatics theoretic techniques. In this chapter, we address the difficulties of mining the huge amount of streaming data and the necessity of interpreting the profiling results in an understandable way.

In Chapter 8, we provide a comprehensive overview of available machine-learning technologies in privacy-preserving data mining. In this chapter, we concentrate on how data-mining techniques lead to privacy breach and how privacy-preserving data mining achieves data protection via machine-learning methods. Privacy-preserving data mining is a new area, and we hope to inspire research beyond the foundations of data mining and privacy-preserving data mining.

In Chapter 9, we describe the emerging challenges in fixed computing or mobile applications and existing and potential countermeasures using machine-learning methods in cybersecurity. We also explore how the emerging cyber threats may evolve in the future and what corresponding strategies can combat threats. We describe the emerging issues in network monitoring, profiling, and privacy preservation and the emerging challenges in intrusion detection, especially those challenges for anomaly detection systems.

Authors

Dr. Sumeet Dua is currently an Upchurch endowed associate professor and the coordinator of IT research at Louisiana Tech University, Ruston, Louisiana. He received his PhD in computer science from Louisiana State University, Baton Rouge, Louisiana.

His areas of expertise include data mining, image processing and computational decision support, pattern recognition, data warehousing, biomedical informatics, and heterogeneous distributed data integration. The National Science Foundation (NSF), the National Institutes of Health (NIH), the Air Force Research Laboratory (AFRL), the Air Force Office of Sponsored Research (AFOSR), the National Aeronautics and Space Administration (NASA), and the Louisiana Board of Regents (LA-BoR) have funded his research with over \$2.8 million. He frequently serves as a study section member (expert panelist) for the National Institutes of Health (NIH) and panelist for the National Science Foundation (NSF)/CISE Directorate. Dr. Dua has chaired several conference sessions in the area of data mining and is the program chair for the *Fifth International Conference on Information Systems, Technology, and Management* (ICISTM-2011). He has given more than 26 invited talks on data mining and its applications at international academic and industry arenas, has advised more than 25 graduate theses, and currently advises several graduate students in the discipline. Dr. Dua is a coinventor of two issued U.S. patents, has (co-)authored more than 50 publications and book chapters, and has authored or edited four books. Dr. Dua has received the Engineering and Science Foundation Award for Faculty Excellence (2006) and the Faculty Research Recognition Award (2007), has been recognized as a distinguished researcher (2004–2010) by the Louisiana Biomedical Research Network (NIH-sponsored), and has won the Outstanding Poster Award at the NIH/NCI caBIG—NCRI Informatics Joint Conference; Biomedical Informatics without Borders: From Collaboration to Implementation. Dr. Dua is a senior member of the IEEE Computer Society, a senior member of the ACM, and a member of SPIE and the American Association for Advancement of Science.

Dr. Xian Du is a research associate and postdoctoral fellow at the Louisiana Tech University, Ruston, Louisiana. He worked as a postdoctoral researcher at the Centre National de la Recherche Scientifique (CNRS) in the CREATIS Lab, Lyon, France, from 2007 to 2008 and served as a software engineer in Kikuze Solutions Pte. Ltd., Singapore, in 2006. He received his PhD from the Singapore–MIT Alliance (SMA) Programme at the National University of Singapore in 2006.

Dr. Xian Du's current research focus is on high-performance computing using machine-learning and data-mining technologies, data-mining applications for cybersecurity, software in multiple computer operational environments, and clustering theoretical research. He has broad experience in machine-learning applications in industry and academic research at high-level research institutes. During his work in the CREATIS Lab in France, he developed a 3D smooth active contour technology for knee cartilage MRI image segmentation. He led a small research and development group to develop color control plug-ins for an RGB color printer to connect to the Windows® system through image processing GDI functions for Kikuze Solutions. He helped to build an intelligent e-diagnostics system for reducing mean time to repair wire-bonding machines at National Semiconductor Ltd., Singapore (NSC). During his PhD dissertation research at the SMA, he developed an intelligent color print process control system for color printers. Dr. Du's major research interests are machine-learning and data-mining applications, heterogeneous data integration and visualization, cybersecurity, and clustering theoretical research.

Chapter 1

Introduction

Many of the nation's essential and emergency services, as well as our critical infrastructure, rely on the uninterrupted use of the Internet and the communications systems, data, monitoring, and control systems that comprise our cyber infrastructure. A cyber attack could be debilitating to our highly interdependent Critical Infrastructure and Key Resources (CIKR) and ultimately to our economy and national security.

Homeland Security Council

National Strategy for Homeland Security, 2007

The ubiquity of cyberinfrastructure facilitates beneficial activities through rapid information sharing and utilization, while its vulnerabilities generate opportunities for our adversaries to perform malicious activities within the infrastructure.* Because of these opportunities for malicious activities, nearly every aspect of cyberinfrastructure needs protection (Homeland Security Council, 2007).

Vulnerabilities in cyberinfrastructure can be attacked horizontally or vertically. Hence, cyber threats can be evaluated horizontally from the perspective of the attacker(s) or vertically from the perspective of the victims. First, we look at cyber threats vertically, from the perspective of the victims. A variety of adversarial agents such as nation-states, criminal organizations, terrorists, hackers, and other malicious users can compromise governmental homeland security through networks.

* Cyberinfrastructure consists of digital data, data flows, and the supportive hardware and software. The infrastructure is responsible for data collection, data transformation, traffic flow, data processing, privacy protection, and the supervision, administration, and control of working environments. For example, in our daily activities in cyberspace, we use health Supervisory Control and Data Acquisition (SCADA) systems and the Internet (Chandola et al., 2009).

For example, hackers may utilize personal computers remotely to conspire, proselytize, recruit accomplices, raise funds, and collude during ongoing attacks. Adversarial governments and agencies can launch cyber attacks on the hardware and software of the opponents' cyberinfrastructures by supporting financially and technically malicious network exploitations.

Cyber criminals threaten financial infrastructures, and they could pose threats to national economies if recruited by the adversarial agents or terrorist organizations. Similarly, private organizations, e.g., banks, must protect confidential business or private information from such hackers. For example, the disclosure of business or private financial data to cyber criminals can lead to financial loss via Internet banking and related online resources. In the pharmaceutical industry, disclosure of protected company information can benefit competitors and lead to market-share loss. Individuals must also be vigilant against cyber crimes and malicious use of Internet technology.

As technology has improved, users have become more tech savvy. People communicate and cooperate efficiently through networks, such as the Internet, which are facilitated by the rapid development of digital information technologies, such as personal computers and personal digital assistants (PDAs). Through these digital devices linked by the Internet, hackers also attack personal privacy using a variety of weapons, such as viruses, Trojans, worms, botnet attacks, rootkits, adware, spam, and social engineering platforms.

Next, we look at cyber threats horizontally from the perspective of the victims. We consider any malicious activity in cyberspace as a cyber threat. A cyber threat may result in the loss of or damage to cyber components or physical resources. Most cyber threats are categorized into one of three groups according to the intruder's purpose: stealing confidential information, manipulating the components of cyberinfrastructure, and/or denying the functions of the infrastructure. If we evaluate cyber threats horizontally, we can investigate cyber threats and the subsequent problems. We will focus on intentional cyber crimes and will not address breaches caused by normal users through unintentional operations, such as errors and omissions, since education and proper habits could help to avoid these threats.* We also will not explain cyber threats caused by natural disasters, such as accidental breaches caused by earthquakes, storms, or hurricanes, as these threats happen suddenly and are beyond our control.

1.1 Cybersecurity

To secure cyberinfrastructure against intentional and potentially malicious threats, a growing collaborative effort between cybersecurity professionals and researchers from institutions, private industries, academia, and government agencies has engaged in

* We define a normal cyber user as an individual or group of individuals who do not intend to intrude on the cybersecurity of other individuals.

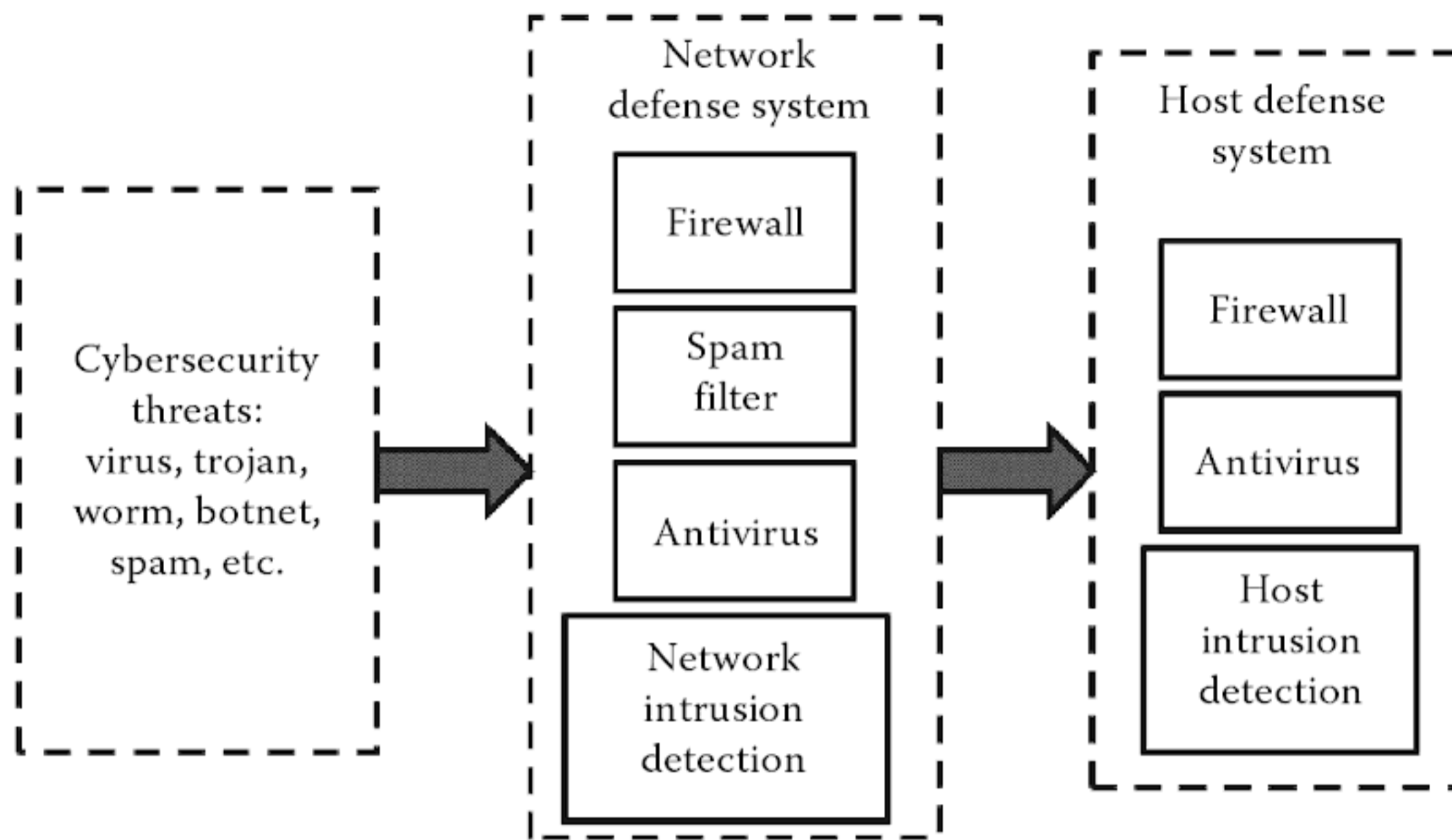


Figure 1.1 Conventional cybersecurity system.

exploiting and designing a variety of cyber defense systems. Cybersecurity researchers and designers aim to maintain the confidentiality, integrity, and availability of information and information management systems through various cyber defense systems that protect computers and networks from hackers who may want to intrude on a system or steal financial, medical, or other identity-based information.*

As shown in Figure 1.1, conventional cybersecurity systems address various cybersecurity threats, including viruses, Trojans, worms, spam, and botnets. These cybersecurity systems combat cybersecurity threats at two levels and provide network- and host-based defenses. Network-based defense systems control network flow by network firewall, spam filter, antivirus, and network intrusion detection techniques. Host-based defense systems control upcoming data in a workstation by firewall, antivirus, and intrusion detection techniques installed in hosts.

Conventional approaches to cyber defense are mechanisms designed in firewalls, authentication tools, and network servers that monitor, track, and block viruses and other malicious cyber attacks. For example, the Microsoft Windows® operating system has a built-in Kerberos cryptography system that protects user information. Antivirus software is designed and installed in personal computers and cyberinfrastructures to ensure customer information is not used maliciously. These approaches create a protective shield for cyberinfrastructure.

However, the vulnerabilities of these methods are ubiquitous in applications because of the flawed design and implementation of software and network

* The three requirements of cybersecurity correspond to the three types of intentional threats: confidentiality signifies the ability to prevent sensitive data from being disclosed to third parties; integrity ensures the infrastructure is complete and accurate, and availability refers to the accessibility of the normal operations of cyberinfrastructures, such as delivering and storing data.

infrastructure. Patches have been developed to protect the cyber systems, but attackers continuously exploit newly discovered flaws. Because of the constantly evolving cyber threats, building defense systems for discovered attacks is not enough to protect users. Higher-level methodologies are also required to discover the embedded and lurking cyber intrusions and cyber intrusion techniques, so that a more reliable security cyberinfrastructure can be utilized.

Many higher-level adaptive cyber defense systems can be partitioned into components as shown in Figure 1.2. Figure 1.2 outlines the five-step process for those defense systems. We discuss each step below.

Data-capturing tools, such as Libpcap for Linux[®], Solaris BSM for SUN[®], and Winpcap for Windows[®], capture events from the audit trails of resource information sources (e.g., network). Events can be host-based or network-based depending on where they originate. If an event originates with log files, then it is categorized as a host-based event. If it originates with network traffic, then it is categorized as a network-based event. A host-based event includes a sequence of commands executed by a user and a sequence of system calls launched by an application, e.g., send mail. A network-based event includes network traffic data, e.g., a sequence of internet protocol (IP) or transmission control protocol (TCP) network packets. The data-preprocessing module filters out the attacks for which good signatures have been learned.

A feature extractor derives basic features that are useful in event analysis engines, including a sequence of system calls, start time, duration of a network flow, source IP and source port, destination IP and destination port, protocol,

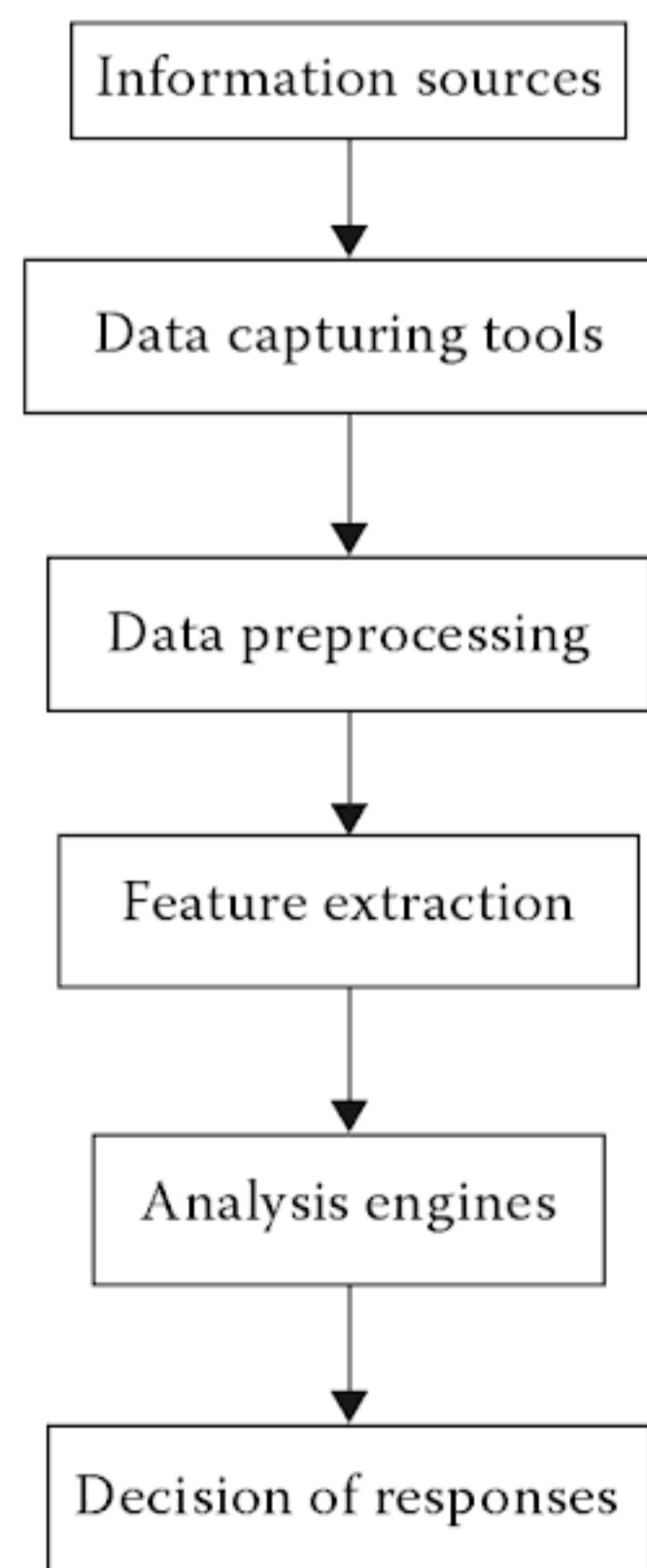


Figure 1.2 Adaptive defense system for cybersecurity.

number of bytes, and number of packets. In an analysis engine, various intrusion detection methods are implemented to investigate the behavior of the cyberinfrastructure, which may or may not have appeared before in the record, e.g., to detect anomalous traffic. The decision of responses is deployed once a cyber attack is identified. As shown in Figure 1.2, analysis engines are the core technologies for the generation of the adaptation ability of the cyber defense system. As discussed above, the solutions to cybersecurity problems include proactive and reactive security solutions.

Proactive approaches anticipate and eliminate vulnerabilities in the cyber system, while remaining prepared to defend effectively and rapidly against attacks. To function correctly, proactive security solutions require user authentication (e.g., user password and biometrics), a system capable of avoiding programming errors, and information protection [e.g., privacy-preserving data mining (PPDM)]. PPDM protects data from being explored by data-mining techniques in cybersecurity applications. We will discuss this technique in detail in Chapter 8. Proactive approaches have been used as the first line of defense against cybersecurity breaches. It is not possible to build a system that has no security vulnerabilities. Vulnerabilities in common security components, such as firewalls, are inevitable due to design and programming errors.

The second line of cyber defense is composed of reactive security solutions, such as intrusion detection systems (IDSs). IDSs detect intrusions based on the information from log files and network flow, so that the extent of damage can be determined, hackers can be tracked down, and similar attacks can be prevented in the future.

1.2 Data Mining

Due to the availability of large amounts of data in cyberinfrastructure and the number of cyber criminals attempting to gain access to the data, data mining, machine learning, statistics, and other interdisciplinary capabilities are needed to address the challenges of cybersecurity. Because IDSs use data mining and machine learning, we will focus on these areas. Data mining is the extraction, or “mining,” of knowledge from a large amount of data. The strong patterns or rules detected by data-mining techniques can be used for the nontrivial prediction of new data. In nontrivial prediction, information that is implicitly presented in the data, but was previously unknown is discovered. Data-mining techniques use statistics, artificial intelligence, and pattern recognition of data in order to group or extract behaviors or entities. Thus, data mining is an interdisciplinary field that employs the use of analysis tools from statistical models, mathematical algorithms, and machine-learning methods to discover previously unknown, valid patterns and relationships in large data sets, which are useful for finding hackers and preserving privacy in cybersecurity.

Data mining is used in many domains, including finance, engineering, biomedicine, and cybersecurity. There are two categories of data-mining methods: supervised and unsupervised. Supervised data-mining techniques predict a hidden function using training data. The training data have pairs of input variables and output labels or classes. The output of the method can predict a class label of the input variables. Examples of supervised mining are classification and prediction. Unsupervised data mining is an attempt to identify hidden patterns from given data without introducing training data (i.e., pairs of input and class labels). Typical examples of unsupervised mining are clustering and associative rule mining.

Data mining is also an integral part of knowledge discovery in databases (KDDs), an iterative process of the nontrivial extraction of information from data and can be applied to developing secure cyberinfrastructures. KDD includes several steps from the collection of raw data to the creation of new knowledge. The iterative process consists of the following steps: data cleaning, data integration, data selection, data transformation, data mining, pattern evaluation, and knowledge representation, as described below.

- Step 1.* During data cleaning, which is also known as data cleansing, noise and irrelevant data are removed from the collection.
- Step 2.* Data integration combines data from multiple and heterogeneous sources into one database.
- Step 3.* Data-selection techniques allow the user to obtain a reduced representation of the data set to keep the integrity of the original data set in a reduced volume.
- Step 4.* In data transformation, the selected data is transformed into suitable formats.
- Step 5.* Data mining is the stage in which analysis tools are applied to discover potentially useful patterns.
- Step 6.* Pattern evaluation identifies interesting and useful patterns using given validation measures.
- Step 7.* In knowledge representation, the final phase of the knowledge-discovery process, discovered knowledge is presented to the users in visual forms.

Data-mining techniques are used to aid in the development of predictive models that enable a real-time cyber response after a sequence of cybersecurity processes, which include real-time data sampling, selection, analysis and query, and mining peta-scale data to classify and detect attacks and intrusions on a computer network (Denning, 1987; Lee and Stolfo, 1998; Axelsson, 2000; Chandola et al., 2006; Homeland Security Council, 2007). Learning user patterns and/or behaviors is critical for intrusion detection and attack predictions. Learning these behaviors is important, as they can identify and describe structural patterns in the data automatically and theoretically explain data and predict patterns. Automatic and theoretic learning require complex computation that calls for abundant machine-learning algorithms. We will discuss the concept of machine learning in Section 1.3.