# APPLICATIONS OF DATA MINING IN COMPUTER SECURITY

# ADVANCES IN INFORMATION SECURITY

## Sushil Jajodia
*Consulting editor*

*Center for Secure Information Systems*
*George Mason University*
*Fairfax, VA 22030-4444*

*email: jajodia@gmu.edu*

*Additional information about this series can be obtained from*
*www.wkap.nl/series.htm/ADIS.*

# APPLICATIONS OF

# DATA MINING IN

# COMPUTER SECURITY

*edited by*

**Daniel Barbará**
**Sushil  Jajodia**
*George Mason University*
*U.S.A.*

# Series Foreword

# ADVANCES IN INFORMATION SECURITY

## Sushil Jajodia
*Consulting Editor*

*Center for Secure Information Systems*
*George Mason University*
*Fairfax, VA 22030-4444*

*email: jajodia@gmu.edu*

Welcome to the sixth volume of the Kluwer International Series on ADVANCES IN INFORMATION SECURITY. The goals of this series are, one, to establish the state of the art of, and set the course for future research in information security and, two, to serve as a central reference source for advanced and timely topics in information security research and development. The scope of this series includes all aspects of computer and network security and related areas such as fault tolerance and software assurance.

ADVANCES IN INFORMATION SECURITY aims to publish thorough and cohesive overviews of specific topics in information security, as well as works that are larger in scope or contain more detailed background information than can be accommodated in shorter survey articles. The series also serves as a forum for topics that may not have reached a level of maturity to warrant a comprehensive textbook treatment.

The success of this series depends on contributions by researchers and developers such as you. If you have an idea for a book that is appropriate for this series, I encourage you to contact me. I would be happy to discuss any potential projects with you. Additional information about this series can be obtained from www.wkap.nl/series.htm/ADIS.

**About This Volume**

This sixth volume of the series is entitled *APPLICATIONS OF DATA MINING IN COMPUTER SECURITY,* edited by Daniel Barbarà and Sushil Jajodia.

Computer intrusions are becoming commonplace and outpacing our capacity to detect, analyze, and counteract them. Since intrusions usually leave traces in the audit data trails, it is only natural to think about this problem in a data-centered way. Some research groups have been successfully using data mining techniques for effectively implementing tools to detect and analyze intrusions.

This volume offers nine articles from leading researchers; eight of these articles focus on the use of data mining for intrusion detection, including one that surveys the state of modern intrusion detection using data mining approaches and another that critically examines these approaches. The last article deals with the application of data mining to computer forensics. Collectively, these articles provide a comprehensive summary of current findings in this fruitful research field.

<div align="right">

SUSHIL JAJODIA
Consulting Editor

</div>

# Contents

7
Adaptive Model Generation 153
*Andrew Honig, Andrew Howard, Eleazar Eskin and Sal Stolfo*

8
Proactive Intrusion Detection 195

# List of Figures

# List of Tables

# Preface

Data mining is becoming a pervasive technology in activities as diverse as using historical data to predict the success of a marketing campaign, looking for patterns in financial transactions to discover illegal activities, or analyzing genome sequences. From this perspective, it was just a matter of time for the discipline to reach the important area of computer security. This book presents a collection of research efforts on the use of data mining in computer security.

Data mining has been loosely defined as the process of extracting information from large amounts of data. In the context of security, the information we are seeking is the knowledge of whether a security breach has been experienced, and, if the answer is yes, who is the perpetrator. This information could be collected in the context of discovering intrusions that aim to breach the privacy of services, or data in a computer system or, alternatively, in the context of discovering evidence left in a computer system as part of a criminal activity.

This collection concentrates heavily on the use of data mining in the area of intrusion detection. The reason for this is twofold. First, the volume of data dealing with both network and host activity is so large that it makes it an ideal candidate for using data mining techniques. Second, intrusion detection is an extremely critical activity. To understand this it is enough to look at the current statistics. Ten major government agencies accounting for 99% of the federal budget have been compromised in the recent past. In the year 2000, a massive, coordinated attack successfully brought down some of the major e-commerce web sites in the United States. Moreover, it is estimated that less than 4% of the attacks are actually detected or reported. As a society, we have become extremely dependent of the use of information systems, so much so that the danger of serious disruption of crucial operations is frightening. As a result, it is no surprise that researchers have produced a relatively large volume of work in the area of data mining in support of intrusion detection.

The rest of the work presented in this volume addresses the application of data mining to an equally pressing area: computer forensics. This area has widened recently to address activities such as law enforcement using digital evidence. Although the amount of work is not as large as in intrusion detection, computer forensics proves to be a fruitful arena for research in data mining techniques.

Data mining holds the promise of being an effective tool to help security activities and, in some sense, the proof of its applicability can be found in the pages of this book. However, there is still a long road to travel and we hope that this volume will inspire researchers and practitioners to undertake some steps in this direction.

## Acknowledgments

DANIEL BARBARÁ

SUSHIL JAJODIA

FAIRFAX, VA