



IRM PRESS

701 E. Chocolate Avenue, Suite 200, Hershey PA 17033-1240, USA
Tel: 717/533-8845; Fax 717/533-8661; URL-<http://www.irm-press.com>

ITB11703

This chapter appears in the book, *Web and Information Security*
edited by Elena Ferrari and Bhavani Thuraisingham © 2006, Idea Group Inc.

Chapter VII

Sanitization and Anonymization of Document Repositories

Yücel Saygin, Sabanci University, Turkey

Dilek Hakkani-Tür, AT&T Labs—Research, USA

Gökhan Tür, AT&T Labs—Research, USA

Abstract

Information security and privacy in the context of the World Wide Web (WWW) are important issues that are still being investigated. However, most of the present research is dealing with access control and authentication-based trust. Especially with the popularity of WWW as one of the largest information sources, privacy of individuals is now as important as the security of information. In this chapter, our focus is text, which is probably the most frequently seen data type in the WWW. Our aim is to highlight the possible threats to privacy that exist due to the availability of document repositories and sophisticated tools to browse

and analyze these documents. We first identify possible threats to privacy in document repositories. We then discuss a measure for privacy in documents with some possible solutions to avoid or, at least, alleviate these threats.

Introduction

Information has been published in various forms throughout the history, and sharing information has been one of the key aspects of development. The Internet revolution and World Wide Web (WWW) made publishing and accessing information much easier than it used to be. However, widespread data collection and publishing efforts on the WWW increased the privacy concerns since most of the gathered data contain private information. Privacy of individuals on the WWW may be jeopardized via search engines and browsers or sophisticated text mining tools that can dig through mountains of Web pages. Privacy concerns need to be addressed since they may hinder data collection efforts and reduce the number of publicly available databases that are extremely important for research purposes such as in machine learning, data mining, information extraction/retrieval, and natural language processing.

In this chapter, we consider the privacy issues that may originate from publishing data on the WWW. Since text is one of the most frequently and conveniently used medium in the WWW to convey information, our main focus will be text documents. We basically tackle the privacy problem in two phases. The first phase, referred to as *sanitization*, aims to protect the privacy of the contents of the text against possible threats. Sanitization basically deals with the automatic identification of named entities such as sensitive terms, phrases, proper names, and numeric values (e.g., credit card numbers) in a given text, and modification of them with the purpose of hiding private information. The second phase, called *anonymization*, makes sure that the classification tools cannot predict the owner or author of the text.

In the following sections, we first provide the taxonomy of possible threats. In addition to that, we propose a privacy metric for document databases based on the notion of k -anonymity together with a discussion of the methods that can be used for preserving privacy.

14 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage: www.igi-global.com/chapter/sanitization-anonymization-document-repositories/31086

Related Content

SEC-CMAC A New Message Authentication Code Based on the Symmetrical Evolutionist Ciphering Algorithm

Bouchra Echandouri, Fouzia Omary, Fatima Ezzahra Ziani and Anas Sadak (2018). *International Journal of Information Security and Privacy* (pp. 16-26).

www.irma-international.org/article/sec-cmac-a-new-message-authentication-code-based-on-the-symmetrical-evolutionist-ciphering-algorithm/208124

Computer Security Practices and Perceptions of the Next Generation of Corporate Computer Users

S. E. Kruck and Faye P. Teer (2011). *Pervasive Information Security and Privacy Developments: Trends and Advancements* (pp. 255-265).

www.irma-international.org/chapter/computer-security-practices-perceptions-next/45815

Cybersecurity Curricular Guidelines

Matt Bishop, Diana Burley and Lynn A. Fletcher (2019). *Cybersecurity Education for Awareness and Compliance* (pp. 158-180).

www.irma-international.org/chapter/cybersecurity-curricular-guidelines/225923

Networking Fundamentals

(2019). *Constructing an Ethical Hacking Knowledge Base for Threat Awareness and Prevention* (pp. 106-118).

www.irma-international.org/chapter/networking-fundamentals/218416

Goals and Practices in Maintaining Information Systems Security

Zippy Erlich and Moshe Zviran (2010). *International Journal of Information Security and Privacy* (pp. 40-50).

www.irma-international.org/article/goals-practices-maintaining-information-systems/50307