

**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

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: <u>www.igi-</u> <u>global.com/chapter/sanitization-anonymization-document-</u>

repositories/31086

### **Related Content**

# A Six-View Perspective Framework for System Security: Issues, Risks, and Requirements

Surya B. Yadav (2010). International Journal of Information Security and Privacy (pp. 61-92).

www.irma-international.org/article/six-view-perspective-framework-system/43057

### Secure Multiparty Computation via Oblivious Polynomial Evaluation

Mert Özararand Attila Özgit (2013). *Theory and Practice of Cryptography Solutions for Secure Information Systems (pp. 253-278).* www.irma-international.org/chapter/secure-multiparty-computation-via-oblivious/76519

### ASKARI: A Crime Text Mining Approach

Caroline Chibelushi, Bernadette Sharpand Hanifa Shah (2008). *Information Security and Ethics: Concepts, Methodologies, Tools, and Applications (pp. 1701-1718).* www.irma-international.org/chapter/askari-crime-text-mining-approach/23187

### Dual Image-Based Dictionary Encoded Data Hiding in Spatial Domain

Giridhar Maji, Sharmistha Mandaland Soumya Sen (2020). *International Journal of Information Security and Privacy (pp. 83-101).* 

www.irma-international.org/article/dual-image-based-dictionary-encoded-data-hiding-in-spatialdomain/247428

#### Social Media and Cyber Security: Investigating the Risk in Nigeria

Desmond Onyemechi Okochaand Damilare J. Agbele (2022). *Handbook of Research on Cyber Law, Data Protection, and Privacy (pp. 50-63).* www.irma-international.org/chapter/social-media-and-cyber-security/300904