# Data-Centric Systems and Applications

Bing Liu

# Web Data Mining

Exploring Hyperlinks,
Contents, and Usage Data

With 177 Figures

Springer

*Bing Liu*

Department of Computer Science
University of Illinois at Chicago
851 S. Morgan Street
Chicago, IL 60607-7053
USA
liub@cs.uic.edu

*To my parents, my wife Yue and children Shelley and Kate*

# Preface

The rapid growth of the Web in the last decade makes it the largest publicly accessible data source in the world. Web mining aims to discover useful information or knowledge from Web hyperlinks, page contents, and usage logs. Based on the primary kinds of data used in the mining process, Web mining tasks can be categorized into three main types: Web structure mining, Web content mining and Web usage mining. Web structure mining discovers knowledge from hyperlinks, which represent the structure of the Web. Web content mining extracts useful information/knowledge from Web page contents. Web usage mining mines user access patterns from usage logs, which record clicks made by every user.

The goal of this book is to present these tasks, and their core mining algorithms. The book is intended to be a text with a comprehensive coverage, and yet, for each topic, sufficient details are given so that readers can gain a reasonably complete knowledge of its algorithms or techniques without referring to any external materials. Four of the chapters, structured data extraction, information integration, opinion mining, and Web usage mining, make this book unique. These topics are not covered by existing books, but yet they are essential to Web data mining. Traditional Web mining topics such as search, crawling and resource discovery, and link analysis are also covered in detail in this book.

Although the book is entitled *Web Data Mining*, it also includes the main topics of data mining and information retrieval since Web mining uses their algorithms and techniques extensively. The data mining part mainly consists of chapters on association rules and sequential patterns, supervised learning (or classification), and unsupervised learning (or clustering), which are the three most important data mining tasks. The advanced topic of partially (semi-) supervised learning is included as well. For information retrieval, its core topics that are crucial to Web mining are described. This book is thus naturally divided into two parts. The first part, which consists of Chaps. 2–5, covers data mining foundations. The second part, which contains Chaps. 6–12, covers Web specific mining.

Two main principles have guided the writing of this book. First, the basic content of the book should be accessible to undergraduate students, and yet there are sufficient in-depth materials for graduate students who plan to

pursue Ph.D. degrees in Web data mining or related areas. Few assumptions are made in the book regarding the prerequisite knowledge of readers. One with a basic understanding of algorithms and probability concepts should have no problem with this book. Second, the book should examine the Web mining technology from a practical point of view. This is important because most Web mining tasks have immediate real-world applications. In the past few years, I was fortunate to have worked directly or indirectly with many researchers and engineers in several search engine and e-commerce companies, and also traditional companies that are interested in exploiting the information on the Web in their businesses. During the process, I gained practical experiences and first-hand knowledge of real-world problems. I try to pass those non-confidential pieces of information and knowledge along in the book. The book, thus, should have a good balance of theory and practice. I hope that it will not only be a learning text for students, but also a valuable source of information/knowledge and even ideas for Web mining researchers and practitioners.

## Acknowledgements

# Table of Contents

# Part I:  Data Mining Foundations

# Part II:   Web Mining