

Chapman & Hall/CRC
Data Mining and Knowledge Discovery Series

Intelligent Technologies for Web Applications

Priti Srinivas Sajja
Rajendra Akerkar



CRC Press
Taylor & Francis Group

A CHAPMAN & HALL BOOK

Intelligent Technologies for Web Applications

Chapman & Hall/CRC

Data Mining and Knowledge Discovery Series

SERIES EDITOR

Vipin Kumar

University of Minnesota

Department of Computer Science and Engineering

Minneapolis, Minnesota, U.S.A.

AIMS AND SCOPE

This series aims to capture new developments and applications in data mining and knowledge discovery, while summarizing the computational tools and techniques useful in data analysis. This series encourages the integration of mathematical, statistical, and computational methods and techniques through the publication of a broad range of textbooks, reference works, and handbooks. The inclusion of concrete examples and applications is highly encouraged. The scope of the series includes, but is not limited to, titles in the areas of data mining and knowledge discovery methods and applications, modeling, algorithms, theory and foundations, data and knowledge visualization, data mining systems and tools, and privacy and security issues.

PUBLISHED TITLES

UNDERSTANDING COMPLEX DATASETS:

DATA MINING WITH MATRIX DECOMPOSITIONS

David Skillicorn

COMPUTATIONAL METHODS OF FEATURE SELECTION

Huan Liu and Hiroshi Motoda

CONSTRAINED CLUSTERING: ADVANCES IN ALGORITHMS, THEORY, AND APPLICATIONS

Sugato Basu, Ian Davidson, and Kiri L. Wagstaff

KNOWLEDGE DISCOVERY FOR COUNTERTERRORISM AND LAW ENFORCEMENT

David Skillicorn

MULTIMEDIA DATA MINING: A SYSTEMATIC INTRODUCTION TO CONCEPTS AND THEORY

Zhongfei Zhang and Ruofei Zhang

NEXT GENERATION OF DATA MINING

Hillol Kargupta, Jiawei Han, Philip S. Yu, Rajeev Motwani, and Vipin Kumar

DATA MINING FOR DESIGN AND MARKETING

Yukio Ohsawa and Katsutoshi Yada

THE TOP TEN ALGORITHMS IN DATA MINING

Xindong Wu and Vipin Kumar

GEOGRAPHIC DATA MINING AND KNOWLEDGE DISCOVERY, SECOND EDITION

Harvey J. Miller and Jiawei Han

TEXT MINING: CLASSIFICATION, CLUSTERING, AND APPLICATIONS

Ashok N. Srivastava and Mehran Sahami

BIOLOGICAL DATA MINING

Jake Y. Chen and Stefano Lonardi

INFORMATION DISCOVERY ON ELECTRONIC HEALTH RECORDS

Vagelis Hristidis

TEMPORAL DATA MINING

Theophano Mitsa

RELATIONAL DATA CLUSTERING: MODELS, ALGORITHMS, AND APPLICATIONS

Bo Long, Zhongfei Zhang, and Philip S. Yu

KNOWLEDGE DISCOVERY FROM DATA STREAMS

João Gama

STATISTICAL DATA MINING USING SAS APPLICATIONS, SECOND EDITION

George Fernandez

INTRODUCTION TO PRIVACY-PRESERVING DATA PUBLISHING:

CONCEPTS AND TECHNIQUES

Benjamin C. M. Fung, Ke Wang, Ada Wai-Chee Fu, and Philip S. Yu

HANDBOOK OF EDUCATIONAL DATA MINING

Cristóbal Romero, Sebastian Ventura, Mykola Pechenizkiy, and Ryan S.J.d. Baker

DATA MINING WITH R: LEARNING WITH CASE STUDIES

Luís Torgo

MINING SOFTWARE SPECIFICATIONS: METHODOLOGIES AND APPLICATIONS

David Lo, Siau-Cheng Khoo, Jiawei Han, and Chao Liu

DATA CLUSTERING IN C++: AN OBJECT-ORIENTED APPROACH

Guojun Gan

MUSIC DATA MINING

Tao Li, Mitsunori Ogihara, and George Tzanetakis

MACHINE LEARNING AND KNOWLEDGE DISCOVERY FOR

ENGINEERING SYSTEMS HEALTH MANAGEMENT

Ashok N. Srivastava and Jiawei Han

SPECTRAL FEATURE SELECTION FOR DATA MINING

Zheng Alan Zhao and Huan Liu

ADVANCES IN MACHINE LEARNING AND DATA MINING FOR ASTRONOMY

Michael J. Way, Jeffrey D. Scargle, Kamal M. Ali, and Ashok N. Srivastava

FOUNDATIONS OF PREDICTIVE ANALYTICS

James Wu and Stephen Coggeshall

INTELLIGENT TECHNOLOGIES FOR WEB APPLICATIONS

Priti Srinivas Sajja and Rajendra Akerkar

Intelligent Technologies for Web Applications

Priti Srinivas Sajja
Rajendra Akerkar



CRC Press

Taylor & Francis Group

Boca Raton London New York

CRC Press is an imprint of the
Taylor & Francis Group, an **informa** business
A CHAPMAN & HALL BOOK

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

© 2012 by Taylor & Francis Group, LLC
CRC Press is an imprint of Taylor & Francis Group, an Informa business

No claim to original U.S. Government works
Version Date: 20120424

International Standard Book Number-13: 978-1-4398-7164-5 (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 CRC Press Web site at
<http://www.crcpress.com>

Dedication

To Srinivas and Abhignya

Priti Srinivas Sajja

To my parents, Ashalata and Arvind Akerkar

Rajendra Akerkar

Contents

Preface.....	xix
Authors	xxv

Part I: Introduction to the Web, machine learning, new AI techniques, and web intelligence

Chapter 1 Introduction to World Wide Web	3
1.1 Brief history of the Web and the Internet.....	3
1.2 Blogs.....	4
1.3 Tweets.....	6
1.4 Wikis.....	7
1.4.1 Improving wiki content reliability, quality, and security.....	8
1.5 Collaborative mapping	10
1.6 Aggregation technologies.....	10
1.7 Open platforms, application programming interface, and programming tools.....	12
1.8 Web intelligence.....	13
1.9 Intelligence in web applications	15
1.10 Organization of this book	16
Exercises	18
References	19
 Chapter 2 Machine learning concepts.....	 21
2.1 Introduction.....	21
2.2 Linear regression	22
2.3 Supervised learning: Classification	23
2.3.1 Evaluation measures	25
2.3.2 Decision trees	27
2.4 Support vector machines.....	28
2.5 Nearest neighbor classifiers	29
2.6 Unsupervised learning: clustering	29
2.6.1 k -Means clustering.....	30
2.6.2 Difference between clustering and nearest neighbor prediction.....	32
2.6.3 Probabilistic clustering	33
2.7 Hidden Markov models.....	33

2.8	Bayesian methods	34
2.8.1	Bayes theorem	34
2.8.2	Naïve Bayes.....	35
2.8.3	Bayesian networks	36
2.9	Reinforcement learning	36
2.10	Applications of machine learning.....	37
2.10.1	Speech recognition	37
2.10.2	Computer vision	38
2.10.3	Robotics.....	38
2.10.4	Software engineering and programming language	39
2.10.5	Computer games.....	39
2.10.6	Machine learning for web	40
2.11	Conclusion	40
	Exercises	41
	References	41
Chapter 3 Overview of constituents for the new artificial intelligence		43
3.1	Foundations of the new artificial intelligence and knowledge-based system	43
3.1.1	Knowledge-based systems	43
3.1.2	Limitations of symbolic systems	43
3.2	Fuzzy systems	44
3.2.1	Fuzzy set and fuzzy logic.....	44
3.2.2	Fuzzy membership function.....	45
3.2.3	Forms and operations of fuzzy functions.....	46
3.2.4	Fuzzy relations and operations on fuzzy relations	46
3.2.5	Fuzzy rule-based systems	47
3.2.6	Applications of fuzzy logic	47
3.3	Artificial neural networks.....	48
3.3.1	Working of an artificial neuron	49
3.3.2	Architectures of artificial neural network	50
3.3.2.1	Multilayer perceptron architecture	50
3.3.2.2	Kohonen architecture.....	52
3.3.3	Applications of artificial neural network	52
3.4	Genetic algorithms and evolutionary computing.....	53
3.4.1	Basic principles of genetic algorithms	54
3.4.2	An example of genetic algorithm to optimize a function	56
3.4.3	Applications of genetic algorithms	56
3.5	Rough sets.....	57
3.5.1	Applications of rough sets.....	58
3.6	Soft computing.....	59
3.6.1	Applications of soft computing	59
3.7	Benefits of the new AI to World Wide Web	60
	Exercises	60
	References	60
Chapter 4 Web intelligence.....		61
4.1	Internet, web, grid, and cloud	61
4.1.1	Components of typical web	62
4.1.2	Characteristics and benefits of the Web	63

4.2	Introduction to web intelligence	64
4.2.1	Semantic web	66
4.2.2	Social intelligence	68
4.2.3	Search engine techniques	68
4.2.4	Web knowledge management	69
4.2.5	Web information retrieval and filtering	69
4.2.6	Web mining	69
4.2.7	Web agents	69
4.2.8	Human-computer integration	70
4.3	Perspectives of WI	70
4.4	Levels of WI	72
4.4.1	Imparting intelligence at basic infrastructural level	72
4.4.2	Imparting intelligence at knowledge level	72
4.4.3	Imparting intelligence at interface level	72
4.4.4	Imparting intelligence at application level	73
4.5	Goal of WI	73
4.6	Characteristics of web intelligence	73
4.6.1	Openness	73
4.6.2	Intelligent	73
4.6.3	Secured	74
4.6.4	User friendly	74
4.6.5	Agent based	74
4.6.6	Interoperability	74
4.6.7	Global knowledge base	74
4.7	Challenges and issues of WI	75
4.7.1	Nature of knowledge	75
4.7.2	Volume, complexity, and unstructured environment of web	75
4.7.3	Development methods, protocols, security, and quality standards	75
4.7.4	Weak support from AI	75
4.8	Wisdom web	75
4.8.1	Autonomic regulation of functionalities between the resources	76
4.8.2	Embedding knowledge into the Web	76
4.8.3	Improving access mechanisms	76
4.9	Web-based support systems	76
4.10	Designing an intelligent web	77
4.11	Future of WI	77
	Exercises	78
	References	78

Part II: Information retrieval, mining, and extraction of content from the Web

Chapter 5	Web information retrieval	83
5.1	Introduction	83
5.1.1	Managing web data	83
5.1.2	Context and web IR	84
5.2	Typical web search engines	85
5.2.1	Introduction to web crawler	86

5.2.2	Some early work in the area of web crawlers	87
5.2.3	Google searching	89
5.3	Architecture of a web crawler.....	89
5.4	Distributed crawling	92
5.5	Focused spiders/crawlers	92
5.5.1	Architecture of the focused crawler	94
5.5.2	Operational phases for the focused crawler	94
5.5.3	Measuring relevance of the focused crawlers	95
5.6	Collaborative crawling.....	95
5.7	Some tools and open source for web crawling.....	96
5.8	Information retrieval: beyond searching	97
5.9	Models of information retrieval	99
5.9.1	Boolean model and its variations	99
5.9.2	Vector space model.....	99
5.9.3	Probabilistic models	100
5.9.4	Latent semantic indexing	100
5.10	Performance measures in IR	100
5.11	Natural language processing in conjunction with IR	101
5.11.1	Generic NLP architecture of IR	102
5.12	Knowledge-based system for information retrieval.....	103
5.13	Research trends.....	106
5.13.1	Semantic information.....	106
5.13.2	Multimedia data	107
5.13.3	Opinion retrieval	107
5.14	Conclusion	107
	Exercises	107
	References	108
Chapter 6	Web mining.....	111
6.1	Introduction to web mining.....	111
6.1.1	Web as a graph.....	112
6.2	Evolution of web mining techniques.....	112
6.3	Process of web mining.....	113
6.4	Web content mining	115
6.4.1	Classification	115
6.4.2	Cluster analysis.....	116
6.4.3	Association mining	117
6.4.4	Structured data extraction.....	118
6.4.5	Unstructured content extraction	118
6.4.6	Template matching	120
6.5	Web usage mining.....	121
6.5.1	Activities in web usage mining.....	121
6.5.2	Retrieval sources for web usage mining.....	121
6.5.3	Cleaning and data abstraction.....	122
6.5.4	Identification of required information	122
6.5.5	Pattern discovery and analysis.....	123
6.6	Web structure mining.....	123
6.6.1	HITS concept	123
6.6.2	PageRank method.....	124

6.7	Sensor web mining: architecture and applications	125
6.8	Web mining software.....	127
6.9	Opinion mining.....	127
6.9.1	Feature-based opinion mining.....	129
6.10	Other applications using AI for web mining	129
6.11	Future research directions	130
	Exercises	131
	References	131
Chapter 7	Structured data extraction.....	133
7.1	Preliminaries.....	133
7.1.1	Structured data	133
7.1.2	Information extraction.....	135
7.1.3	Evaluation metrics	137
7.1.4	Approaches to information extraction	138
7.1.5	Free, structured, and semistructured text	138
7.1.6	Web documents.....	139
7.2	Wrapper induction	140
7.2.1	Wrappers.....	140
7.2.2	From information extraction to wrapper generation	140
7.2.3	Wrapper generation.....	141
7.2.3.1	Semiautomated wrapper generation.....	142
7.2.3.2	Automated wrapper generation	145
7.2.4	Inductive learning of wrappers	151
7.3	Locating data-rich pages.....	152
7.3.1	Finding tables.....	152
7.3.2	Identifying similarities	152
7.3.3	Heuristics on product properties	153
7.3.4	Human intrusion	153
7.4	Systems for wrapper generation.....	153
7.4.1	Structured and semistructured web pages	154
7.4.1.1	ShopBot.....	154
7.4.1.2	Wrapper induction environment.....	154
7.4.1.3	SoftMealy	155
7.4.1.4	Supervised learning algorithm for inducing extraction rules	156
7.4.1.5	WebMantic	158
7.4.2	Semistructured and unstructured web pages	159
7.4.2.1	Robust automated production of IE rules	160
7.4.2.2	Sequence rules with validation.....	161
7.4.2.3	WHISK.....	162
7.5	Applications and commercial systems.....	163
7.5.1	Examples of applications.....	164
7.5.2	Commercial systems	164
7.5.2.1	Junglee	165
7.5.2.2	Jango	165
7.5.2.3	MySimon	166
7.6	Summary	166
	Exercises	167
	References	167

Part III: Semantic web and web knowledge management

Chapter 8	Semantic web	173
8.1	Introduction to semantic web	173
8.2	Metadata	174
8.2.1	Dublin core metadata standard	176
8.2.2	Metadata objectives	176
8.2.2.1	Simplicity of creation and maintenance	176
8.2.2.2	Commonly understood semantics	179
8.2.2.3	International scope	179
8.2.2.4	Extensibility	179
8.2.2.5	Interoperability	179
8.3	Layered architecture of semantic web	180
8.3.1	Unicode and uniform resource identifier	180
8.3.2	Extensible markup language	180
8.3.3	Resource description framework	180
8.3.4	RDF schema	181
8.3.5	Ontology	181
8.3.6	Logic and proof	181
8.3.7	Trust	181
8.4	Refined architecture of semantic web	181
8.5	Ontology and ontology constructs	182
8.5.1	Extensible markup language	184
8.5.2	Resource description framework	187
8.5.3	Web ontology language	188
8.5.3.1	OWL full	190
8.5.3.2	OWL DL	190
8.5.3.3	OWL lite	190
8.5.4	Ontology interchange language	191
8.5.5	OWL2 profile	191
8.5.6	SPARQL	194
8.5.6.1	Result syntaxes	196
8.5.6.2	Query for relationships	196
8.5.6.3	Transform data with CONSTRUCT	196
8.5.6.4	OPTIONAL	196
8.5.6.5	Negation	197
8.6	Meta-ontology	198
8.7	Ontology tools and editors	198
8.8	Annotation tools	199
8.9	Inference engines	199
8.10	Semantic web applications	200
8.10.1	Search engine	200
8.10.2	Semantic web portals	201
8.10.3	Catalog management and thesaurus	201
8.10.4	Call center	201
8.10.5	e-Learning	201
8.10.6	Tourism	202
8.10.7	Publishing	204
8.10.8	Community and social projects	204

8.10.9	e-Commerce	205
8.10.10	Health care.....	205
8.10.11	Digital heritage.....	205
8.10.12	Open archives	206
8.11	Semantic web interoperability and web mining	207
8.12	Semantic web and social communities	207
8.13	Semantic web and intelligent search	208
8.14	Semantic web research issues.....	210
	Exercises	210
	References	211

Chapter 9 Web knowledge management213

9.1	About knowledge	213
9.2	Knowledge management fundamentals.....	213
9.2.1	Architecture of the knowledge management process.....	215
9.2.2	Benefits of knowledge management.....	216
9.2.3	Challenges of knowledge management	217
9.3	Ontology revisited	217
9.3.1	Ontology examples.....	217
9.3.2	Ontology classification.....	219
9.3.3	Parameters to build ontology.....	220
9.3.4	Standards and interoperability for ontology	221
9.3.5	Ontology on the Web	221
9.4	Utilization of knowledge management methodologies on semantic web	222
9.4.1	Literature review	223
9.4.2	General architecture for web knowledge management.....	224
9.4.3	Semantically enhanced knowledge management	225
9.4.3.1	Semantic wiki	226
9.4.3.2	Semantic annotation tools	227
9.4.4	Issues and challenges.....	228
9.5	Exchanging knowledge in virtual entities.....	228
9.5.1	Virtual world	228
9.5.2	Virtual organizations	229
9.5.3	Knowledge management and intelligent techniques within virtual entities	230
9.5.4	Virtual communities and semantic web	232
9.6	Case study.....	233
9.7	Building the World Wide Why	234
9.8	Conclusion and applications	234
	Exercises	235
	References	235

Chapter 10 Social network intelligence239

10.1	Introduction to social networking	239
10.1.1	Web patterns and social ecosystem	242
10.1.2	Types of social networks.....	242
10.2	Friend-of-a-friend	243
10.3	Semantically interlinked online communities	248

10.4	Social network analysis	249
10.5	Social network data	251
10.6	hCard and XFN	252
10.7	Advantages and disadvantages of social networking	254
10.7.1	Advantages of social networking	254
10.8	Social graph application programming interface	254
10.9	Social search and artificial intelligence	255
10.9.1	Intelligent social networks	257
10.10	Research future	257
	Exercises	258
	References	258

Part IV: Agent-based web, security issues, and human–computer interaction

Chapter 11	Agent-based web	263
11.1	Introduction	263
11.2	Agents	264
11.2.1	Characteristics and advantages	264
11.2.2	Agents and objects	266
11.2.3	Agents and web services	266
11.3	Typology of agents	267
11.3.1	Collaborative agent	267
11.3.2	Interface agent	269
11.3.3	Mobile agent	269
11.3.4	Information agent	269
11.3.5	Intelligent agent	270
11.3.6	Hybrid agent	272
11.4	Multiagent systems	272
11.4.1	Multiagent system framework	273
11.4.2	Communication between agents	274
11.5	Agent-based web	274
11.5.1	Generic architecture of agent-based web	275
11.5.2	Example agents	277
11.5.2.1	Agent for query and information retrieval	277
11.5.2.2	Filtering agent	277
11.5.2.3	Interface agent	278
11.5.2.4	Personal assistance agent	278
11.5.2.5	e-Commerce agent	279
11.5.2.6	e-Communities and agents	279
11.5.2.7	Ontology management	280
11.6	Hybridization of mobile agent and interface agent: A case for personalized content representation	281
11.6.1	Mobile agents: Characteristics and working	281
11.6.2	Hybridization of a mobile agent with an interface agent	282
11.6.3	Personalized content representation through the hybrid agent	284
11.7	Case study	286
11.7.1	Multiagent system for oil company	286
11.7.2	RETSINA calendar agent	287

11.7.2.1	OpenStudy.com	289
11.7.2.2	Cobot.....	290
11.8	Conclusion	291
	Exercises	291
	References	292
Chapter 12	Web security	293
12.1	Introduction.....	293
12.2	Web vulnerabilities.....	294
12.2.1	Scripting languages.....	294
12.2.2	Understanding communication	295
12.2.3	Injection flaws	296
12.2.4	Cross-site scripting.....	297
12.2.5	Cross-site request forgery.....	298
12.2.6	Phishing attacks.....	298
12.2.7	Information leakage	299
12.2.8	Browsers compromising privacy.....	299
12.3	Web server protection	299
12.3.1	Firewall.....	299
12.3.2	Intrusion detection system	300
12.4	Security and privacy	301
12.5	Contributions of AI for security issues	302
	Exercises	304
	References	305
Chapter 13	Human–web interactions	307
13.1	Introduction.....	307
13.2	Features of a good website	308
13.2.1	Content	308
13.2.2	Information organization	308
13.2.3	Performance.....	309
13.2.4	Compatibility.....	309
13.2.5	Visual design	309
13.2.6	Interaction design	309
13.3	What is interaction?	309
13.3.1	Common interaction styles	310
13.3.2	Three-dimensional interactions	310
13.4	Interaction design and related parameters	311
13.5	Usability	313
13.5.1	World usability day	314
13.6	Process of interaction design	314
13.6.1	Know your users, their requirements, and identify the objective of the system.....	315
13.6.2	Check the feasibility and cost–benefit ratio of different alternatives	315
13.6.3	Build selected alternatives considering content.....	315
13.6.4	Test the developed interaction design	315
13.6.5	Conduct postimplementation review and update if required	315
13.7	Conceptual models of interaction	315
13.8	Interface.....	316

13.9	Interface design methods	317
13.9.1	Activity-centered design.....	317
13.9.2	Body storming.....	317
13.9.3	Contextual design.....	318
13.9.4	Focus group	318
13.9.5	Iterative design.....	318
13.9.6	Participatory design	319
13.9.7	Task analysis.....	319
13.9.8	User-centered design.....	319
13.9.9	Usage-centered design.....	320
13.9.10	User scenario	320
13.9.11	Value-sensitive design.....	320
13.9.12	Wizard of Oz experiment	320
13.10	Tools for human–web interaction.....	321
13.11	Interaction evaluation methods.....	321
13.11.1	Cognitive walk-through	321
13.11.2	Heuristic evaluation	321
13.11.3	Review-based evaluation.....	322
13.11.4	Evaluating through user participation	322
13.11.5	Evaluating implementations	322
13.12	Human–computer interaction and human–web interaction	322
13.13	Issues in human–web interactions.....	323
13.14	Support of AI for human–web interactions	323
13.14.1	Searching, retrieval and filtering, and semantic search	324
13.14.2	Native language interface and fuzzy logic for web applications	325
13.14.3	Knowledge management and knowledge representation on the Web.....	325
13.14.4	Agent-based systems.....	326
13.14.5	Modeling users experience and usability for better web interactions	326
13.14.6	Intelligent web mining	326
13.14.7	Interacting with smart environments.....	326
13.15	Case studies.....	327
13.15.1	MIT intelligent room	327
13.15.2	MediaBlocks system	327
13.15.3	PhotoHelix	328
13.16	Research applications.....	328
13.17	Conclusion	329
	Exercises	329
	References	330

Preface

The Web is becoming the largest data repository in the world and presents a key driving force for a large spectrum of information technology (IT). To develop effective and intelligent web applications and services, it is critical to discover useful knowledge through analyzing large amounts of content, hidden content structures, or usage patterns of web data resources. To achieve such a goal, a variety of techniques in diverse research areas need to be integrated, including natural language processing, information extraction, information retrieval, information filtering, knowledge representation, knowledge management, machine learning, databases, data mining, web mining, text mining, agent, human–computer interaction, and the semantic web. These integrated techniques should address the key challenges from the heterogeneous and dynamic nature of web contents and usage patterns.

Within the past ten years, the Web research community has brought to maturity a comprehensive set of foundational technology components, both at the conceptual level and in the form of prototypes and software.

Intended readers

This book describes the basics as well as the latest trends in the area of an integrated approach instead of an edited volume of papers/chapters. The book provides a detailed review of issues for web researchers. With extensive use of examples and more than 100 illustrations, as well as bibliographical notes, end-of-chapter exercises, and glossaries, to clarify complex material and demonstrate practical applications, this book can serve as a senior undergraduate-level book. It can also serve as a good reference for researchers and practitioners who deal with the various problems involving semantics, intelligent techniques for web ontologies, and the semantic web.

Understanding web-related concepts, studying the underlying standards and technical components, and putting all of this together into concrete terms require a substantial amount of effort. This book provides comprehensive and easy-to-follow coverage on both the “what-is” and “how-to” aspects of web-related technologies.

In particular, this book is written keeping the following readers in mind:

- Software engineers and developers who are interested in learning the intelligent and semantic web technology in general
- Web application developers who are interested in studying the intelligent web technologies and in constructing web applications
- Researchers who are interested in the research and development of intelligent and semantic technologies

- Undergraduate and graduate students in computer science departments, whose main area of focus is the intelligent and semantic web
- Practitioners in related engineering fields

The prerequisites needed to understand the concepts in this book include the following:

- Working knowledge of a programming language
- Basic understanding of the Web, including its main technical components such as URL, HTML, and XML

Salient features

This book has the following salient features:

- Makes all fundamental as well as in-depth material available at one place in an integrated manner
- Provides a more concrete organization than an edited volume
- Incorporates new topics on artificial intelligence (AI), thus making the book more effective and helpful in solving problems
- Integrates illustrations and examples to support pedagogical exposition
- Equips the reader with the necessary information in order to obtain hands-on experience of the topics of discussion
- Facilitates experimentation of the content discussed in the book by making available fundamental tools, research directions, practice questions, and additional reading material
- Integrates all material, yet allows each chapter to be used or studied independently
- Supplies further tools and information at the associated website for instructors and students

Outline of the chapters

The book is organized into four parts. Part I provides an introduction to the Web, machine learning, new AI techniques, and web intelligence.

Chapter 1 describes introductory concepts such as a brief history of the Web and the Internet. It also discusses the latest trends on the Web such as blogs, tweets, wikis, etc. Collaborative mapping, aggregation technologies, open platforms, tools, and application programming interfaces (APIs) are discussed in this chapter. The chapter also describes the organization of the content.

Chapter 2 reviews machine learning that has made its way from AI into web applications and technologies. It presents the capabilities of machine learning methods and provides ideas on how these methods could be useful for web intelligence. The chapter establishes fundamentals such as linear regression, estimation, generalization, supervised learning, unsupervised learning, reinforcement learning, hidden Markov models, and Bayesian networks.

Chapter 3 covers the new AI and knowledge-based system (KBS) and discusses the limitations of the typical symbolic AI and the need of bio-inspired AI for the Web. The most essential and widely employed material pertaining to neural networks, genetic algorithms, fuzzy systems, and rough sets are discussed in brief with their possible advantages.

Chapter 4 explores the basic roles as well as practical impacts of artificial intelligence and advanced information technology for the next generation of web-based systems, services, and environments. The chapter also presents the concept of wisdom web.

Part II is dedicated to information retrieval, mining, and extraction of content from the Web.

Web information retrieval is another important aspect linked to web intelligence. Web spiders, distributed spiders, focused spiders, search engine mechanisms, personalized search techniques, and natural language processing (NLP) in conjunction with effective retrieval are discussed in Chapter 5. This chapter also presents architectures of knowledge-based systems for information retrieval from the Web.

Web mining is the application of machine learning (especially data mining) techniques to web-based data for the purpose of learning or extracting knowledge. Web mining methodologies can generally be classified into one of three distinct categories: web usage mining, web structure mining, and web content mining. Chapter 6 discusses these methodologies along with suitable applications.

Chapter 7 introduces the concept of information extraction to facilitate structured data extraction. Information extraction (IE) is a technology enabling relevant content to be extracted from textual information available electronically. It plays a crucial role for researchers and professionals as well as for other end users who have to deal with vast amounts of information from the Internet. This chapter focuses on wrapper induction as well as semiautomatic and automatic wrapper generation along with a suitable case study.

Part III is dedicated to the semantic web and web knowledge management.

Chapter 8 establishes the semantic web as an immediate extension of the Web in which the meaning (semantics) of content and services on the Web is defined along with the content. Embedding of such semantics makes it possible for the Web to “understand” the content and satisfy the requests of people and machines to use the Web. The chapter discusses metadata, metadata standards, layered architecture of semantic web, and tools and ontology constructs such as resource description framework (RDF), web ontology language (OWL), and extensible markup language (XML). Ontology spectrum, meta-ontology, editors, inference and annotation tools, etc., are also included. It also discusses web applications such as semantic search, social communities, and semantic web research issues.

The Web encompasses a large amount of content organized heterogeneously. For effective retrieval and better access of the content available on the Web, it is necessary to use suitable knowledge representation, knowledge use, and knowledge-sharing techniques. Chapter 9 discusses various knowledge management techniques for the Web. It also suggests a generic architecture on the top of the semantic web for knowledge management.

Chapter 10 combines the concepts and the methods of two fields, namely, the semantic web and social networks, which, together, aid in the analysis of the social web and in the design of a new class of applications that combine human intelligence with machine processing. The chapter presents the application of semantic web technologies to the social web that forms a network of interlinked and semantically enabled content and knowledge. It also provides readers with an understanding of the key concepts and methods of both the fields and describes a simple real-world application incorporating social and semantic metadata.

Part IV discusses additional topics such as agent-based web, security issues, and human-computer interaction.

An agent is an entity that is autonomous, independent, and cooperative. It does intended work on behalf of the user. To carry out various web activities and support web

functionalities in a structured manner, one may take the help of agents. Chapter 11 discusses agent typology, intelligent agents, agents for the Web, web services, and case studies. Considering the technologies discussed within the aforementioned chapters, some agents can be designed to fit into the framework of a multi-agent web. One such possible framework of a multi-agent system is discussed in this chapter. The chapter also elaborates on applications suitable for the framework suggested.

Chapter 12 discusses issues related to web security. It reviews different AI and machine learning methods concerning security, privacy, and reliability issues of cyberspace. It also enables readers to discover the types of methods at their disposal, summarizes the state of the practice in this important area, and provides a classification of existing work. The topics include security management and governance, network security and authentication, intrusion detection, trust management, access control, and privacy.

The expectations from the Web are ever increasing, and the Web will also evolve accordingly. However, the facilities offered by such a giant organization would be made more effective with better interface. Chapter 13 focuses on human–web interactions. It defines web interaction and identifies interaction applications. Topics such as interactive information search/retrieval, interactive query expansion, personalization, user profiling, visualization, user interfaces, usability, web adaptation, and interactive authoring/annotation for the semantic web are discussed in this chapter along with other similar applications.

Use as a book

The book can be covered in a total of approximately 40–45 lecture hours (plus 20–30 hours dedicated to exercises and hands-on practice).

Parts I and II can be covered as a complete course in about 30 taught hours. Such a course requires a significant amount of additional practical activity, normally consisting of several exercises from each chapter and a project involving the design and implementation of a web application.

Parts III and IV can be covered in a second course. They can alternatively be integrated in part within an extended first course. In advanced, project-centered courses, the study of current technology can be accompanied by a project dedicated to the development of technological components. The advanced course can be associated with further readings or with a research-oriented seminar series.

Acknowledgments

The organization and the contents of this book have benefited from our experience in teaching the subject in various contexts. All the students attending those courses, dispersed over many schools and countries (Sardar Patel University, the International School of Information Management, Saint Mary's University, American University of Armenia, and SIBER-India), deserve our deepest gratitude. Some of these students have class-tested rough drafts and incomplete notes and have contributed to their development, improvement, and correction. Similarly, we would like to thank the staff from IT companies and government organizations who attended our courses for professionals and helped us learn the practical aspects that we have tried to convey in this book. We would also like to thank all the colleagues who have contributed, directly or indirectly, to the development of this book, through discussions on course organization or the actual revision of drafts. They include Pawan Lingras, Terje Aaberge, Svein Ølnes, David Camacho, Henry Hexmoor, and Darshan Choksi.

We thank the reviewers of this edition for a number of very useful suggestions concerning the organization of the book and the specific content of chapters.

We also thank Aastha Sharma, David Fausel, Sarah Morris, the staff at CRC Press, and Remya Divakaran (SPi Global) who have contributed to the birth of this book.

Finally, we express our gratitude to our families for their love, support, and patience during the preparation of the book. We also thank our families for reminding us that there are things in life beyond writing books.

Priti Srinivas Sajja

Rajendra Akerkar

Authors

Priti Srinivas Sajja joined the faculty of the Department of Computer Science, Sardar Patel University, Gujarat, India, in 1994 and is presently working as an associate professor. She received her MS (1993) and PhD (2000) in computer science from Sardar Patel University. Her research interests include knowledge-based systems, soft computing, multi-agent systems, and software engineering. She has more than 100 publications in books, book chapters, journals, and in the proceedings of national and international conferences. Three of her publications have won best research paper awards. Dr. Sajja is the coauthor of *Knowledge-Based Systems*. She supervises the work of seven doctoral research students. She is also the principal investigator of a major research project funded by the University Grants Commission (UGC), India. She serves as a member on the editorial boards of many international science journals and has served as a program committee member for various international conferences.

Rajendra Akerkar is a professor/senior researcher at Western Norway Research Institute (Vestlandsforskning), Norway. His research and teaching experience includes over 20 years in academia, spanning different universities in Asia, Europe, and North America. As the founder of Technomathematics Research Foundation (TMRF), he is instrumental in ensuring that the organization lends a platform for research in India. Under his leadership, TMRF has become a well-known organization among the research community worldwide.

Akerkar's current research agenda focuses on learning and language—how each works in the human and how they can be replicated in a machine. He received a DAAD fellowship in 1990 and was also awarded the prestigious BOYSCASTS Young Scientist Award of the Department of Science & Technology, Government of India, in 1997. He is the editor in chief of the *International Journal of Computer Science & Applications* and an associate editor of the *International Journal of Metadata, Semantics, and Ontologies*. Akerkar serves as a member of the scientific committees of several international conferences and also serves on the editorial boards of international journals in computer science. He has authored 12 books, more than 90 research papers, and has edited 5 volumes of international conferences. He initiated the International Conference Series on Web Intelligence, Mining and Semantics (WIMS). Akerkar has been actively involved in many industrial research and development projects for more than 14 years.

part one

*Introduction to the Web, machine
learning, new AI techniques,
and web intelligence*

Introduction to World Wide Web

1.1 Brief history of the Web and the Internet

The World Wide Web, abbreviated as WWW and commonly known as the Web, has been weaving a variety of solutions for different problems and meeting information requirements of a global audience. It is a system of interlinked hypertext documents in multimedia accessed via Internet, which is defined as network of networks. The dream was conceived by Tim Berners-Lee, who is now director of the World Wide Web Consortium and extending the dream project further in a form of semantic web by adding semantics to the existing web. The Web is developed to be a pool of information to allow collaborators from remote sites to share their ideas and information.

During the year 1980, Tim Berners-Lee built ENQUIRE as a personal database of people using hypertext and software utilities to access the database. The objective was to share data globally without common machines and presentation software. He implemented this system on a newly acquired NeXT workstation. After the invention of supporting hypertext transfer protocol (HTTP) and a web browser named World Wide Web, the first web server and page were created that described the project itself. It was further modified to be used on any machine rather than NeXT. On August 6, 1991, Berners-Lee posted a short summary of the World Wide Web project on the alt.hypertext newsgroup. This date also marked the debut of the Web as a publicly available service on the Internet. According to the summary, the World Wide Web (WWW) project aimed to allow all links to be made to any information anywhere. He invited high-energy physicists and other experts to share data, news, and documentation. Inspired from the message, university-based scientific departments and physics laboratories adopted the concept developed such as Fermilab (Fermi National Accelerator Laboratory for high-energy physics, Batavia, IL) and SLAC (Stanford Linear Accelerator Center, Stanford University, Menlo Park, CA).

There was still no graphical browser available for computers besides the NeXT. This gap was filled in 1992 with the release of Erwise (an application developed at Helsinki University of Technology, Finland) and ViolaWWW (created by Pei-Yuan Wei, which included advanced features such as embedded graphics, scripting, and animation). This gave rise to the development of different web browsers. Some prominent early browsers are Mosaic (now Netscape Navigator) and Cello. Immediately after that, in 1994, the World Wide Web Consortium was founded at Massachusetts Institute of Technology (MIT) with the support of Defense Advance Research Project Agency (DARPA) and European Commission. Berners-Lee made the Web available freely, with no patent and no royalties due. By the end of 1994, the total number of websites has increased, however, the increase is minute (in comparison with) the present standards of 15 million index pages approximately. However, by 1996, the usage of the Web was no longer optional. Earlier, people had identified the possibilities of free publishing and instant worldwide information, but, at present, the Web has opened up the possibility of direct web-based commerce and instantaneous group communications worldwide. The innovation of protocols, standards, and utilities like search engine and e-mail made the Web ubiquitous and fall within the reach

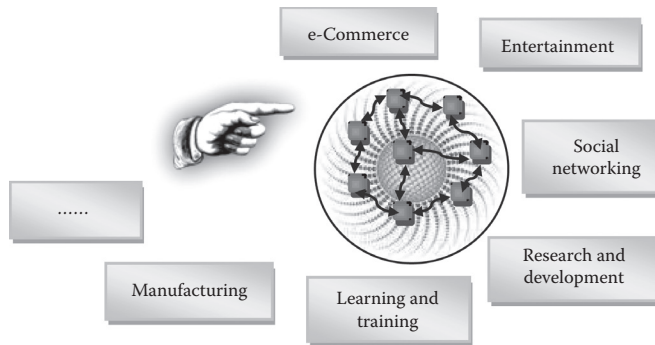


Figure 1.1 Ubiquitous applications of the Web.

of a common man. New ideas such as blogs and social networking are also welcomed in recent times. Some typical application areas of the Web include e-commerce, decision support in various businesses, e-learning, social networking, training, information repository services, manufacturing, and research and development. Figure 1.1 represents ubiquitous application of the Web. Some of the recent innovations are discussed in the next section.

1.2 Blogs

A blog, derived from the term web log (or weblog), is a type of website that is usually maintained by an individual to share details regarding personal views, events, or multimedia materials such as graphics or videos. The term weblog was coined by Jorn Barger on December 17, 1997, to describe the list of links on his Robot Wisdom website that “logged” his Internet wanderings. The short form “blog” was coined by Peter Merholz, who jokingly broke the word *weblog* into the phrase *we blog* in the sidebar of his blog Peterme.com in April–May, 1999. Blog entries are commonly displayed in reverse-chronological order, from most recent entries to the oldest. Blog can also be used as a verb, meaning to maintain or add content to a blog. Many blogs are interactive and allow visitors to leave comments and messages. Such instructiveness makes the blog different from other static websites. Further, a typical website encompasses plenty of web pages, each may have some subpages under a home page following a well-organized approach. A blog is normally a single page of entries by a single author in the most personalized way. There may be archives of older entries optionally to manage the large stuff in a well-organized way. A blog is normally public and accessible to anyone without any formality such as registration. However, to post a personal comment that may require to be processed (replied) further, there is a need of user information such as an e-mail ID. Figure 1.2 represents a typical blog.

Early blogs were simply manually updated components of common websites. However, the evolution of tools to facilitate the production and maintenance of blog is much easier and less technical. Blogs can be hosted by dedicated blog hosting services or on regular web-hosting services. Alternatively, a blog can be developed and maintained using blog software (also called as blogware) like Wordpress (which is a third-party software, see <http://wordpress.org/>) and Blogsmith (<http://www.blogsmith.com/>). The blog software supports the authoring, editing, and publishing of blog posts and comments, with special functions for image management, moderation of posts, and comments. All blog software support authoring, editing, and publishing of entries in the following format:



Figure 1.2 Example of a blog.

- Title or headline of the post
- Body, main content of the post
- Reference or link to other article
- Date and time on which the blog is published
- Blog entries can optionally include comments and categories (or tags)

Generally blogs publish focused content on a particular topic such as home design, sports, mobile systems technology, etc. However, there may be some blogs that provide a variety of content and links to plenty of other locations. Most of the blogs have a few things in common. These include

- Heading and main content
- Archives
- Facility for the readers to comment and contact
- Some useful related links (blogrolls)
- Feeds like really simple syndication (RSS), resource description framework (RDF), or atoms
- Excerpts (summary) and plugins for readymade additional functionalities

A blog may have features like trackbacks and pingbacks in order to allow users to comment on blog posts and link to the posts and comment on and recommend them further. Track back is a way to provide notification between different websites. If a person finds the specific blog content interesting, he may send trackback *ping* to the Web author. Pingback offers the advanced facility of automatically notifying the author that the other person referred the author's post. Pingbacks can be automated and do not send content, whereas trackbacks are manual and send an excerpt of the initiating post. The facility of trackback and pingback aim to provide some control on comments on blog content and hence help in extending authority to blog commenting. Some tools offer a feature of

comment moderation to monitor and control the comments on the different article posts. Here authors (blog publishers) are given rights to manage the comment spam, delete harsh comments, and approve cool comments. Finally, it must be remembered that blogs require continuous, regular, and meaningful update.

1.3 Tweets

A tweet is a small message, post, or status update on some network. The size of a tweet is comparatively smaller than the typical e-mail. The “tweet” originally means a sound of a bird or whistle, which sweetly says something. Tweet is also considered as a real-time microblogging service. The term tweet became popular by a social networking website called Twitter (Twitter Inc., San Francisco, CA) created in 2006 by Jack Dorsey who is an American software architect and businessman. This site offers facility to send and receive tweets up to 140 characters. That is why it is sometimes described as the short message service (SMS) of the Internet. Generally, tweets are displayed on the user’s home page or profile page of the website and are visible by default free of cost. However, it may be restricted to the user’s friends list. Other users may choose to read and subscribe (opt) for some specific users tweets. These subscribers are known as followers. The process of subscribing a user’s tweet is known as “following.” Alternatively, such tweets can be followed on compatible external applications such as smart phones or by SMS available. It is reported that Twitter currently has more than 175 million users. The original project code name for the Twitter service was twttr, inspired by Flickr and the five-character length of American SMS short codes. Flickr is an image-hosting and video-hosting website for online community created by Ludicorp and later acquired by Yahoo! It is a popular website to share and embed personal photographs and images. The basic advantage of using the facility to tweet a message or an item is that one can “post once” and the service will redistribute the item to multiple followers.

Twitter’s Application Programming Interface (API) is based on the representational state transfer (REST) architecture. REST is a collection of network design principles and guidelines that define resources and ways to address and access data. With the REST architecture, the Twitter works with most web syndication formats that help in gathering information from one source and sends it out to various destinations. Twitter is compatible with two of them—RSS and atom syndication format (atom). Both formats retrieve data from one resource and send it to another. A web page administrator can embed RSS/atom code into the code of his or her site. Visitors can subscribe to the syndication service—called feed—and receive an update every time the administrator updates the web page. Twitter uses this feature to allow members to post messages to a network of other Twitter members. In effect, Twitter members subscribe to other members’ feeds. By allowing third-party developers partial access to its API, Twitter allows them to create programs that incorporate Twitter’s services.

Current third-party applications include the following applications:

- *Twitterlicious* and *Twitterific*, two applications that allow users to access Twitter through desktop applications on PCs and MACs, respectively
- *OutTwit*, a Windows application that allows users to access Twitter through the Outlook e-mail program
- *Tweet Scan*, which allows users to search public Twitter posts in real time
- *Twessenger*, which integrates with the Windows Live Messenger 8.1 instant messenger program

- *Twittervision*, which integrates a Twitter feed into Google Maps
- *Flotzam*, which integrates Twitter with Facebook, Flickr, and blogs
- *iTunes to Twitter*, an application that broadcasts the title of the song currently playing in the user's iTunes to his or her network
- *TwitterBox*, a Twitter application that works inside the virtual community of Second Life

The ability to send and receive multiple tweets without any limit and sophisticated interface made twitter very popular in societies and industries. Tweeting is a tool for accessing opinions, decisions, and market information. Tweeting can be used as a very good communication and business tool. It can be used in the field of education for classroom community and collaborative writing within and across schools and institutes.

1.4 Wikis

A Wiki is a collaborative web platform that allows multiple users from different locations to add information at a centralized place in an interactive fashion. The interface for interaction is created in such a way that users do not require any training. The wiki concept was first introduced by Ward Cunningham as "the simplest online database that could possibly work" in the year 1995. This is also considered as writable web or open editing concept. Wiki looks like a simple website with an edit facility/link. Readers of the page can modify it or add content directly through this link. Wikis can be used as a centralized repository for many applications and provides efficient document management. Wiki can be used as a website, as a knowledge base for FAQ system, and as an information sharing utility. As many authors contribute and evaluate information for a specific topic, wiki can enhance the quality and quantity of information. Table 1.1 shows some typical applications that can be benefited by wikis.

On request, wiki shares available information and allows editing of the information through any browser. It also invites users to create new pages and promotes a specified topic. Wiki can create links to new empty pages that do not yet exist to invite users to share their opinion. There are many different ways in which users can edit the wiki content. Figure 1.3 gives an idea how users can work with wiki.

The content of wiki is normally specified with a simplified markup language, sometimes known as "*wiki text*." Some wikis allow dedicated content such as a repository of images or collections of audio. Common facilities that a wiki typically provides is creating new pages, editing content, navigation, and searching. Sometimes external search engine facility can be embedded within a wiki. Majority of wikis provide limited access to hypertext markup language (HTML) and cascade style sheet (CSS). This is the prime reason that generally wikis look simple. Many modern wikis are making "What You See Is What You Get" (WYSIWYG) editing available to users through scripts and controls

Table 1.1 Applications Benefited by Wikis

Project management and up-to-date project status information
Tutorials and e-learning
Internal notice board and institutional news circulation
FAQ management
Online document management and information sharing
Managing groups and social interactions

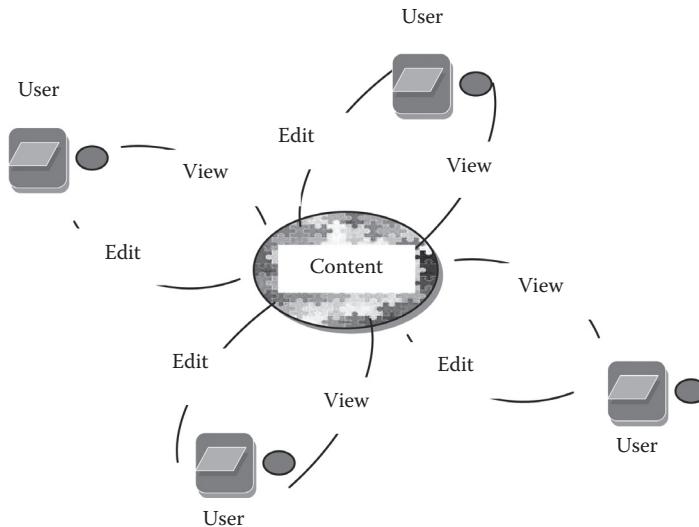


Figure 1.3 Working with wikis.

(like JavaScript or ActiveX controls). The scripts and controls translate the formatting instructions into appropriate wiki text. The translation is transparent and hidden from users to avoid complexity.

Most wikis maintain a log of changes made to wiki pages and maintain versions of the pages. The administrator or author may revert to the older version. With the help of such a revert facility, authors can revert to an older version of the page. Some wikis like MediaWiki insist contributors to provide separate edit summary when they edit a page. The edit summary is an important parameter that justifies the editing and helps in taking the decision whether the page is to be reverted back or not.

For implementation of such collaborative and sharable platform, wiki software is used. The wiki software is implemented on a server and can be executed on different web servers. The content is stored in file systems and changes are stored in a relational database for version management. Some software that are helpful in creating wikis are given in Table 1.2.

Wiki can be developed and used as a standalone system. Such wikis are known as personal wikis. The example of such a personal wiki is the Wikidpad.

1.4.1 Improving wiki content reliability, quality, and security

In spite of its simplicity and advantages, wiki is not considered a reliable source of information. As anybody can contribute in such a system, the content of the system is easily tampered. However, there are genuine contributors and editors, who catch such malicious tampering and correct the mistakes and destructive content for the benefit of others. Many times, users contribute the content that they thought correct and genuine, such as ideas, research information, and experiments. It is not advisable to change such information; however, in this situation, one or more versions showing suggestions and changes can be prepared and a link can be provided adjacent to this. This strategy is also applicable if contributor to wiki does not have rights to edit the content. To improve the reliability, the content pages are ranked by administrator, readers, and contributors. There are

Table 1.2 List of Software for Wiki

ConcourseConnect is a freely available J2EE application with social networking, online community, business directory, and customer relationship management capabilities. This tool supports features like wiki, blog, document management, reviews, online advertising, and project management modules

DokuWiki is a simple-to-use Wiki to facilitate documentation management needs within a small institution. It uses plain text files and has a simple but powerful syntax, which ensures the data files remain readable outside the Wiki

MediaWiki is a popular free web-based wiki software application. It is developed by, and it runs, all the projects of the Wikimedia Foundation, including Wikipedia, Wiktionary, and Wikinews. It is written in the hypertext processor (PHP) programming language and uses a backend database

TiddlyWiki is a HTML/JavaScript-based wiki in which the entire site/wiki is contained in a single file. This tool does not require any server support

Wikidpad is a freeware (open source) personal-use wiki with native support of international characters (Unicode). This software is executable on a single machine. It helps in the implementation of features like storing thoughts, ideas, to-do lists, contacts, and other notes with wiki-like linking between pages

Windows SharePoint Server 2010 has built-in Wiki support. It is built on ASP.Net, C#, and Microsoft SQL Server

Wikia (formerly known as Wikicities) is a free web-hosting service for wikis

sophisticated software programs that help in improving reliability, security, and quality of wiki content by checking and imposing constraints on content and contributors. Table 1.3 suggests some ideas to improve trustworthiness and quality of wiki content.

Wikis that allow only registered users are known as closed wikis. In general, wiki allows anonymous editing by just recording IP address of the machine from which editing is done. In closed wikis, only registered users, whose information (like name and biography) is formally recorded, can edit the wiki content. With this strategy, wiki content is reliable and controllable; however, growth of content on wiki is slower. The popular Wikipedia is an open wiki that allows anonymous editing and records only IP addresses. The countermeasures shown in Table 1.3 help in preventing attacks from malicious contributors, vandalism content, harmful code, and bugs. Besides these, an edit war may occur

Table 1.3 Improving Reliability, Quality, and Security of Wiki Content

Improvement	Countermeasures
Improving reliability of content	By correcting the content manually with the help of experts and administrators By providing links to the edited version of the content By accreditation of users as well as content
Improving quality	By frequent editing, ranking, and filtering the content
Improving security	By allowing only registered users to edit the content By imposing requirement of additional waiting period before allowing any edit By dedicated software support (such as JavaScript) to automatically find vandalism, bugs, and harmful content By preparing a list of malicious sites, if any, within the content pointing to any of the sites from the list