

Data Analytics

Series editors

Longbing Cao, Advanced Analytics Institute, University of Technology, Sydney,
Broadway, NSW, Australia

Philip S. Yu, University of Illinois at Chicago, Chicago, IL, USA

Aims and Goals:

Building and promoting the field of data science and analytics in terms of publishing work on theoretical foundations, algorithms and models, evaluation and experiments, applications and systems, case studies, and applied analytics in specific domains or on specific issues.

Specific Topics:

This series encourages proposals on cutting-edge science, technology and best practices in the following topics (but not limited to):

Data analytics, data science, knowledge discovery, machine learning, big data, statistical and mathematical methods for data and applied analytics,

New scientific findings and progress ranging from data capture, creation, storage, search, sharing, analysis, and visualization,

Integration methods, best practices and typical examples across heterogeneous, interdependent complex resources and modals for real-time decision-making, collaboration, and value creation.

More information about this series at <http://www.springer.com/series/15063>

Chuan Shi · Philip S. Yu

Heterogeneous Information Network Analysis and Applications

 Springer

Chuan Shi
Beijing University of Posts and
Telecommunications
Beijing
China

Philip S. Yu
University of Illinois at Chicago
Chicago, IL
USA

Data Analytics

ISBN 978-3-319-56211-7

ISBN 978-3-319-56212-4 (eBook)

DOI 10.1007/978-3-319-56212-4

Library of Congress Control Number: 2017936890

© Springer International Publishing AG 2017

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

The publisher, the authors and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, express or implied, with respect to the material contained herein or for any errors or omissions that may have been made. The publisher remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Printed on acid-free paper

This Springer imprint is published by Springer Nature

The registered company is Springer International Publishing AG

The registered company address is: Gewerbestrasse 11, 6330 Cham, Switzerland

Preface

The interacting and multi-typed components in the real-world environment constitute interconnected networks, which can be called information networks. These ubiquitous information networks form a critical component of modern information infrastructure. In recent years, the information network analysis has gained extremely wide attentions from researchers in many disciplines, such as computer science, social science, physics. Particularly, the information network analysis has become a mainstream direction in data mining, database and information retrieval fields in the past decades. The basic paradigm is to mine hidden patterns through mining linkage relations from networked data. The information network analysis is also related to the works in social network analysis, link mining, graph mining and network science.

Contemporary information network analyses are usually based on homogeneous information networks, where there is only one type of objects or links in the network. An example is the author collaboration network which only contains the author object and the co-author relation. These homogeneous information networks usually are the simplification of real interacting systems by simply ignoring the heterogeneity of objects and links or only considering one type of links among one type of objects. However, most real interacting systems contain multi-typed interacting components which can be modeled as heterogeneous information networks which include different types of objects and links. For example, the bibliographic database, like DBLP, can be organized as a heterogeneous information network which includes multiple types of objects (e.g., papers, authors, and venues) and links (e.g., written by/writing relations between papers and authors, published/publishing relations between papers and venues). Obviously, the author collaboration network is implicitly contained in the heterogeneous information network, which can be derived from the written by/writing relation between papers and authors.

Compared to homogeneous information network, the heterogeneous information network can effectively fuse more information and contain richer semantics in objects and links, and thus it forms a new development of data mining. Since the concept of heterogeneous information network is first proposed in 2009, it rapidly

became a hot research topic in data mining, and many innovative data mining tasks have been exploited in this kind of networks. In addition, some unique analysis techniques (e.g., meta-path-based mining) are developed to demonstrate the benefits of heterogeneous information networks. Particularly, with the arrival of the era of big data, heterogeneous information networks offer the potential to be an effective way to model and analyze complex objects and their relations in big data.

This book first provides a comprehensive survey of current developments of heterogeneous information network analysis, as well as some novel data mining tasks in this field. This book includes two parts. In the first part, it deeply and comprehensively summarizes the newest developments of this field in Chaps. 1, 2, and 9. This book introduces in-depth understanding of heterogeneous information network in Chap. 1 and investigates the research developments in most data mining tasks in Chap. 2. Furthermore, based on the newest developments and trends, we point out the future research directions in Chap. 9. In the second part, it illustrates the traits of heterogeneous information network analysis through several data mining tasks in Chaps. 3–8. This book presents relevance measure in Chap. 3, ranking and clustering in Chap. 4, recommendation in Chap. 5, fusion learning in Chap. 6, and schema-rich heterogeneous network mining in Chap. 7. Moreover, some interesting prototype systems are discussed in Chap. 8.

The readers of this book are engineers and researchers in the field of data mining, especially social network analysis. It is also suitable for engineers and researchers in artificial intelligences and informatics. More broadly, readers also include those who are interesting in social network analysis in other disciplines, such as statistics, social sciences, physical, and biology. This book can be used in those courses, such as data mining, social network analysis, complex network, advanced artificial intelligences. These courses are suitable for advanced undergraduates or graduate students specializing in computer sciences and related fields. The readers are suggested to quickly understand this field through the first part and deeply study data mining tasks with the second part.

We would like to express our sincere thanks to all those who work with us on this project. First of all, we appreciate Dr. Jiawei Zhang for his contribution in Chap. 6, which makes this book more integrated. Then, we are grateful to our co-authors in the work of heterogeneous information network. They are Xiangnan Kong, Yizhou Sun, Bin Wu, Yitong Li, Zhiqiang Zhang, Jian Liu, Ran Wang, Yuyan Zheng, Jing Zheng, Xiaohuan Cao, Jiawei Hu, Xiaofeng Meng, Chong Zhou, et al. We also wish to thank supporters during writing this book. They are Xin Wan, Xiaoji Chen, Yugang Ji, Houye Ji, Yiding Zhang, Yang Xiao, Binbin Hu, Xiaotian Han, Pudi Chen, Li Song, Govardhana K., Melissa Fearon, Jennifer Malat, et al. In addition, this work is supported by the National Key Basic Research and Department (973) Program of China (No. 2013CB329600), the National Natural Science Foundation of China (No. 61375058 and 61672313), and US National Science Foundation through grant III-1526499. We also thank the supports of these grants. Finally, we thank our families for their wholehearted support throughout this project.

Contents

1	Introduction	1
1.1	Basic Concepts and Definitions	1
1.2	Comparisons with Related Concepts	5
1.3	Example Datasets of Heterogeneous Information Networks	6
1.4	Why Heterogeneous Information Network Analysis	8
	References	9
2	Survey of Current Developments	13
2.1	Similarity Search	13
2.2	Clustering	15
2.3	Classification	17
2.4	Ranking	18
2.5	Link Prediction	19
2.6	Recommendation	21
2.7	Information Fusion	22
2.8	Other Applications	24
	References	24
3	Relevance Measure of Heterogeneous Objects	31
3.1	HeteSim: A Uniform and Symmetric Relevance Measure	31
3.1.1	Overview	31
3.1.2	The HeteSim Measure	33
3.1.3	Experiments	40
3.1.4	Quick Computation Strategies and Experiments	47
3.2	Extension of HeteSim	52
3.2.1	Overview	52
3.2.2	AvgSim: A New Relevance Measure	53
3.2.3	Parallelization of AvgSim	55
3.2.4	Experiments	56

3.3	Conclusion	59
	References	60
4	Path-Based Ranking and Clustering	61
4.1	Meta Path-Based Ranking	61
4.1.1	Overview	61
4.1.2	The HRank Method	63
4.1.3	Experiments	70
4.2	Ranking-Based Clustering	80
4.2.1	Overview	80
4.2.2	Problem Formulation	82
4.2.3	The HeProjI Algorithm	84
4.2.4	Experiments	91
4.3	Conclusions	95
	References	95
5	Recommendation with Heterogeneous Information	97
5.1	Recommendation Based on Semantic Path	97
5.1.1	Overview	97
5.1.2	Heterogeneous Network Framework for Recommendation	99
5.1.3	The SemRec Solution	103
5.1.4	Experiments	108
5.2	Recommendation Based on Matrix Factorization	117
5.2.1	Overview	117
5.2.2	The SimMF Method	118
5.2.3	Experiments	122
5.3	Social Recommendation with Heterogeneous Information	130
5.3.1	Overview	130
5.3.2	The DSR Method	131
5.3.3	Experiments	136
5.4	Conclusions	139
	References	140
6	Fusion Learning on Heterogeneous Social Networks	143
6.1	Network Alignment	143
6.1.1	Overview	143
6.1.2	Terminology Definition and Social Meta Path	144
6.1.3	Cross-Network Network Alignment	148
6.1.4	Experiments	150
6.2	Link Transfer Across Aligned Networks	154
6.2.1	Overview	154
6.2.2	Cross-Network Link Prediction	155
6.2.3	Experiments	159

- 6.3 Synergistic Network Community Detection 166
 - 6.3.1 Overview 166
 - 6.3.2 Cross-Network Community Detection 166
 - 6.3.3 Experiments 170
- 6.4 Conclusions 177
- References 178
- 7 Schema-Rich Heterogeneous Network Mining 181**
 - 7.1 Link Prediction in Schema-Rich Heterogeneous Network 181
 - 7.1.1 Overview 181
 - 7.1.2 The LiPaP Method 183
 - 7.1.3 Experiments 187
 - 7.2 Entity Set Expansion with Meta Path in Knowledge Graph 190
 - 7.2.1 Overview 190
 - 7.2.2 The MP_ESE Method 191
 - 7.2.3 Experiments 195
 - 7.3 Conclusions 198
 - References 198
- 8 Prototype System Based on Heterogeneous Network 201**
 - 8.1 Semantic Recommender System 201
 - 8.1.1 Overview 201
 - 8.1.2 System Architecture 203
 - 8.1.3 System Implementation 204
 - 8.1.4 System Demonstration 208
 - 8.2 Explainable Recommender System 208
 - 8.2.1 Overview 208
 - 8.2.2 Heterogeneous Network-Based Recommendation 210
 - 8.2.3 System Framework 212
 - 8.2.4 System Demonstration 212
 - 8.3 Other Prototype Systems on Heterogeneous Network 215
 - 8.4 Conclusions 216
 - References 216
- 9 Future Research Directions 219**
 - 9.1 More Complex Network Construction 219
 - 9.2 More Powerful Mining Methods 220
 - 9.2.1 Network Structure 221
 - 9.2.2 Semantic Mining 221
 - 9.3 Bigger Networked Data 224
 - 9.4 More Applications 224
 - References 225